Development of a High-Latitude Convection Model by Application of Machine Learning to SuperDARN observations

William A. Bristow¹, Charles Topliff², and Morris B. Cohen²

¹Pennsylvania State University ²Georgia Institute of Technology

November 21, 2022

Abstract

A new model of high-latitude convection derived using machine learning (ML) is presented. The ML algorithm random forests regression was applied to a database of velocity observations from the Super Dual Auroral Radar Network (SuperDARN). The features used to train the model were the IMF components B_x , B_y , and B_z ; the solar wind velocity, v_{sw} ; the auroral indicies, A_u and A_{l} ; and the geomagnetic index, SYM-H. The SuperDARN velocities were separated into north-south, and east-west components and sorted into a magnetic local time - magnetic latitude grid that ran from 55° to the magnetic pole with a bin size of 2° in latitude, and 1-hour in MLT. Separate models were created for each velocity component in each bin of the grid. It is found that even though the models in each bin are independent of one another a coherent convection pattern is formed when the models are viewed in aggregate. The resulting convection pattern responds to changes in the auroral indicies by expanding and contracting in a way that is consistent with expectations for a substorm cycle. Further it is found that the mean-squared difference between predictions of the model and observed values of the velocity are substantially lower than the same quantity calculated for an existing climatology that was not formed with ML techniques

Development of a High-Latitude Convection Model by Application of Machine Learning to SuperDARN observations

3

4

5

W. A. Bristow 1 , C. Topliff 2 , and M. B. Cohen 2

 $^1{\rm Pennsylvania}$ State University $^2{\rm Georgia}$ Institute of Technology

6 Abstract

A new model of high-latitude convection derived using machine learning (ML) is presented. The ML 7 algorithm random forests regression was applied to a database of velocity observations from the Su-8 per Dual Auroral Radar Network (SuperDARN). The features used to train the model were the IMF 9 components B_x , B_y , and B_z ; the solar wind velocity, v_{sw} ; the auroral indicies, A_u and A_l ; and the ge-10 omagnetic index, SYM-H. The SuperDARN velocities were separated into north-south, and east-west 11 components and sorted into a magnetic local time - magnetic latitude grid that ran from 55° to the 12 magnetic pole with a bin size of 2° in latitude, and 1-hour in MLT. Separate models were created for 13 each velocity component in each bin of the grid. It is found that even though the models in each bin 14 are independent of one another a coherent convection pattern is formed when the models are viewed 15 in aggregate. The resulting convection pattern responds to changes in the auroral indicies by expand-16 ing and contracting in a way that is consistent with expectations for a substorm cycle. Further it is found 17 that the mean-squared difference between predictions of the model and observed values of the veloc-18 ity are substantially lower than the same quantity calculated for an existing climatology that was not 19 formed with ML techniques. 20

²¹ **1** Introduction

Climatological convection modeling has been carried out for decades by binning various localized mea-22 surements of the ionospheric plasma velocity or electric field collected over the high-latitude regions 23 versus some set of parameters, most often the interplanetary magnetic field (IMF) components. Such 24 models often are used to drive circulation models such as the Thermospere lonosphere Electrodynamic 25 General Circulation Model (TIEGCM) (Roble & Ridley, 1994) and the Global Ionosphere-Thermosphere 26 Model (GITM) (Ridley et al., 2006). They are also used to constrain the data driven convection pat-27 terns produced from SuperDARN data (Ruohoniemi & Baker, 1998). The measurements have been col-28 lected using a variety of instruments such as satellite based drift meters (e.g. Heelis et al., 1982) or elec-29 tric field booms (e.g. Heppner & Maynard, 1987), incoherent-scatter radar (Foster, 1983), ground-based 30 magnetometers (e.g. Papitashvili et al., 1994), and coherent-scatter radars (e.g. Ruohoniemi & Green-31 wald, 1996). Construction of the models typically has involved grouping observations based upon pre-32 vailing IMF conditions and perhaps some other parameter such as the planetary K-index (k_p) (Heppner 33 & Maynard, 1987), or the geomagnetic Auroral Electrojet Index (A_e) (Weimer, 2005), or the dipole 34

-2-

tilt angle (Thomas & Shepherd, 2018), and then using the binned observations to constrain an expan-

³⁶ sion of the electrostatic potential in a set of orthogonal functions.

The underlying assumption of such a binning is that when repeated, a given set of driving conditions 37 will on average produce the same unique convection pattern. In a general sense, physical reasoning and 38 observations show this to be true. For example, when the IMF is southward, there is magnetic merg-39 ing on the dayside magnetopause and the near-noon field lines connecting from the ionosphere out to 40 the magnetopause are directly influenced by the electric field across the merging region. Those field 41 lines are convected anti-sunward across the polar cap from noon to midnight where they eventually re-42 connect to the field lines from the opposite hemisphere. Once reconnected, the demand for magnetic 43 flux on the dayside causes the field lines to return, following a path through the magnetosphere that 44 maps to the ionosphere just equatorward of the polar cap boundary. The total potential drop across 45 the polar cap should be equal to the projection of the solar wind electric field along the magnetopause 46 reconnection region multiplied by the length of that region. Hence, for a set of solar wind/IMF param-47 eters, the electric field would be the same and if the merging line is of the same length, the potential 48 imposed in the ionosphere would be the same. This scenario leads to a two-celled pattern with flow from 49 noon to midnight across the polar cap and return flow in the opposite direction at lower latitudes. 50

Numerous studies have examined the influence of the driving conditions on various aspects of convec-51 tion. The 1987 article titled "Empirical High-Laittude Electric Field Models" by J. P. Heppner and N. 52 C. Maynard, provides an excellent summary of the patterns that have been observed and how the IMF 53 influences them. In particular, they highlighted the influence of the sign of the IMF y-component on 54 the location and direction of the flow in the dayside throat. Their study contrasts with most others in 55 that rather than binning the observations on a grid and then constraining a functional expansion of the 56 potential with the binned observations, they examined individual satellite passes and categorized them 57 as signatures or "quasi-signatures" and then sorted them based on the IMF. Their result was a set of 58 three basic patterns that covered the majority of southward IMF situations. Those patterns illustrated 59 sharper features (Harang Discontinuity, dayside throat) than are evident in most other models. In ad-60 dition, they examined the influence of k_p and A_e , but only by comparing the average total cross po-61 lar cap potential drop for ranges of the parameters. 62

As illustrated by Hepner and Maynard's discussion of k_p and A_e , there several factors that influence convection that are not accounted for by the instantaneous state of the IMF and solar wind alone. For

-3-

example, while dayside merging converts closed magnetospheric field lines into open polar-cap field lines, 65 night-side merging in the magnetotail converts open field lines to closed, moving them from the po-66 lar cap to the magnetosphere. The diameter of the polar cap and hence the latitude of the convection 67 reversal boundary is determined by the total open magnetic flux, which is determined by the balance 68 between the dayside and nightside merging rates (Siscoe & Huang, 1985). That nightside merging rate 69 is highly variable and depends on the internal state of the magnetosphere. In the growth phase of a 70 substorm the nightside merging rate may be substantially lower than the dayside rate, leading to an 71 expanding polar cap and expanding convection pattern. During a substorm expansion phase the op-72 posite is true and rapid night-side merging can lead to a contracting polar cap and convection pattern. 73 Further, features like the enhancement of the Harang Discontinuity during growth phase (Bristow & 74 Jensen, 2007) change the shape of the pattern in addition to its diameter. 75

To account for some of the dependence on conditions beyond the solar wind and IMF, a new convec-76 tion model was constructed using parameters that provide some indication of the state of the magne-77 tosphere in addition to the solar wind and IMF parameters. Specifically, the auroral indicies A_u and 78 A_l , and the mid-latitude geomagnetic index Sym-H. The auroral indicies give an indication of the 79 strength of convection (A_u) and of the level of substorm activity (A_l) , while the Sym-H index gives 80 an indication of the strength of the ring current, which has been shown to influence the diameter of 81 the auroral oval (Schulz, 1997). These indicies are readily available for use from the NASA OMNI database 82 (King & Papitashvili, 2005), which also provides the solar wind and IMF parameters aligned in time 83 to reflect solar wind propagation delays from the point of observation to the Earth's bow shock. In-84 cluding the magnetospheric parameters increases the dimension of the parameter space to seven, which 85 is fairly large for traditional method of binning the observations. In addition, as will be demonstrated 86 the dependence of the convection velocities on some of the parameters is nonlinear. Because of these 87 two factors, machine learning (ML) was used to form the model. 88

The paper begins with discussion of the SuperDARN data and how the database influenced the choice of algorithms for generating the model. The form of the data motivated producing independent models for the velocity at the points of a latitiude-MLT grid rather than an orthogonal function expansion of the global-scale potential pattern. That discussion is followed by examination the output of the model at a single location and comparison to an existing climatology. Next, the individual models are combined to generate global-scale potential patterns that could be used in the same way as existing clima-

-4-

- ⁹⁵ tologies. Finally, it's demonstrated that the resulting model predicts observations with a lower error than
- 96 the climatology.



97 2 Machine learning implementation

Figure 1: (a) Example convection map from March 21, 2014 created using observations from 0400 UT to 0405 UT, and (b) potential map generated by the Thomas and Shepherd 2018 climatological model for the IMF conditions at the time.

The convection model was based on observations from the Super Dual Auroral Radar Network (Super-98 DARN) (Greenwald et al., 1995) processed with the potential mapping technique, Map-Potential (Ruohoniemi 99 & Baker, 1998). Potential maps were generated for every 5-minute interval from January 1, 2013 to 100 December 31, 2017. The Thomas-Shepherd 2018 (TS18) statistical model (Thomas & Shepherd, 2018) 101 was used as the background for Map-Potential, with the IMF from the OMNI database used to select 102 the patterns. An example pattern from 0400 UT to 0405 UT on March 21, 2014, is shown in Figure 1. 103 The map shows contours of the electrostatic potential along with flow vectors at locations where radar 104 observations constrained the fit. At the time of the plot, the IMF was weekly southward with a neg-105 ative y-component. Figure 1b shows the potential contours predicted by the TS18 model for the con-106

ditions observed at the time from the observed IMF conditions. There are several subtle differences be-107 tween observed pattern and that predicted by the model. First, the total cross-cap potential drop is 108 significantly lower in the observations. The drop is only 30 kV, while the predicted value was 58 kV. The 109 observed pattern is shifted toward midnight and rotated slightly toward dusk. That is, the flow across 110 the polar cap in the observation is from pre-noon to pre-midnight, while the model predicted flow di-111 rectly from noon to midnight. Such differences are typical for any given interval and simply illustrates 112 the inherent variability of the convection pattern and the need for using observations whenever pos-113 sible rather than relying on climatologies. 114

Several methods were considered for generating the ML model. The simplest implementation concep-115 tually would have been to bin the line-of-sight observations from the individual radars similarly to the 116 way they have been used to construct previous models. This works well for constraining global mod-117 els of the potential since the plasma velocity in the F-region ionosphere is very nearly the so called $\mathbf{E} \times \mathbf{B}$ 118 velocity and the observation can be written as the negative gradient of the potential ($\mathbf{E} = -\nabla V$) crossed 119 with the local magnetic field and projected along the line-of-sight. Having an ensemble of observations 120 at different locations is sufficient to constrain the fit for the entire high-latitude region. A second al-121 ternative considered was to form potential patterns and treat them as images. Convolutional neural net-122 works are adept at using such observations, however because the pattern obtained at any given time 123 can be dominated by the influence of the statistical model, any ML model trained on such patterns would 124 tend simply to reproduce the statistical model. 125

Another method considered was to use vectors formed at locations where two or more radars provided 126 measurements. These merged vectors are not influenced by a model, however they would have provided 127 rather sparse coverage and would have required discarding data from locations where only one radar 128 had an observation. In addition, vectors formed in this way can be noisy because the LOS observations 129 from the two radars can potentially be separated in time by tens of seconds which often means they 130 are observing significantly different conditions. Small changes in the azimuth of a flow can result in sig-131 nificant changes in the LOS projections, which are amplified when they are recombined to form a vec-132 tor, especially when the viewing angle between the lines of sight is small. 133

The method we chose was to use the vectors from the SuperDARN potential mapping obtained at locations where there were there were one or more observations contributing to the fit. While the vectors are influenced by the statistical model, having the observation at a location significantly lessens

-6-

that influence. In addition the since the fitted vectors are consistent with an electrostatic potential, which 137 is a strong physical constraint, using them minimizes the impact of noise and radar-to-radar inconsis-138 tencies in the observations. The method had the added benefit of decreasing the size of the database 139 from what would have been required if we had used the observations from the individual radars. One 140 disadvantage of using the data in this way is that values at any given location are not continuous in 141 time, which limited the ML algorithms that could be applied. Individual radars do not have continu-142 ous observations and the longitudinal distribution of radars Is not uniform so some longitude ranges are 143 not covered. The Russian sector is the most obvious illustration of the gap in coverage. 144

The fitted vectors were binned on a magnetic latitude-local-time (MLT) grid from 55° magnetic latitude to the pole with a cell size of 2° in latitude and 1-hour in MLT. The vectors were resolved into north-south and east-west components and written to comma-separated-value files along with the associated value of the IMF vector (B_x , B_y , and B_z), the solar wind velocity (v_{sw}), the Auroral Electrojet Indicies (A_u , A_l), and the Sym-H index, all from the OMNI database.

Figure 2 shows the components of the database for a representative 3-day period for the bin at 65° magnetic latitude and 2000 MLT. As discussed above the data are not continuous in time, which limits the type of ML techniques that can be employed. While it might be desirable to use something like the Long-Short-Term Memory technique (LSTM) since it is a good technique for predicting time evolution based on time-series drivers, continuous observations are necessary to train such a model.

While not continuous, the database is relatively large. Figure 3a shows the number of data points in 155 each grid cell. The highest values illustrated in the figure are in excess of 1-million, however at the low-156 est latitudes some cells have less than 100,000 observations. These locations are often equatorward of 157 the convection zone so there are no usable observations available. The low-latitude extent of signifi-158 cant convection is highly variable in time, like most aspects of convection, and is determined by the 159 magnetospheric drivers and state. In general convection is confined to latitudes above a low-latitude 160 convection boundary referred to as the Heppner-Maynard Boundary (HMB) (Shepherd & Ruohoniemi, 161 2000). Fortunately, by examining the observations from the entire SuperDARN network at a given time 162 it is possible to identify this boundary with some confidence. With that determination, it is possible 163 to assign a zero value at those times to the velocity in cells that lie at latitudes below the HMB. Fig-164 ure 3b shows the density of points including the assignment of zero velocity when a bin is at a lower 165 latitude than the HMB. With this assignment, there are in excess of 400,000 points in all bins between 166

-7-



Figure 2: Time series of observations and feature values in the grid cell at 65° magnetic latitude and 2000 MLT for the interval March 20 - 22, 2014.

¹⁶⁷ 55° and 80°. Above about 85° there are not many observations, so uncertainties will be larger there than
 ¹⁶⁸ in other regions.

Figure 4 illustrates the relationships between the velocity components in the grid cell at 67° magnetic 169 latitude and 1800 MLT, and some of the parameters from the database. In each frame of the figure, 170 the vertical axis is one of the velocity components $(v_{ns} \text{ or } v_{ew})$ and the horizontal axis is one of the 171 database parameters. Pixel color indicates the density of points in a bin. Solid black dots indicate the 172 average velocity in each of the parameter bins. No color scale is provided since the goal is to exam-173 ine trends and not to extract quantitative information. The purpose of examining the data in this way 174 is to select parameters for inclusion as features for training the ML model. If the velocities were un-175 correlated with any of the parameters it would be possible to exclude them and decrease the complex-176 ity of the model. With that in mind, it is still interesting to examine the trends that the plots show. 177 As would be expected for the auroral zone latitude dusk MLT location, the magnitude of the north-178 south component (v_{ns}) is significantly smaller than the east-west component (v_{ew}) . v_{ew} shows a strong 179 dependence on each of the selected parameters, though the dependence is clearly nonlinear for A_u and 180 A_l . Frame 4a illustrates that v_{ew} is negative (westward) for the vast majority of the data, indicating 181



(a) Number of observations

(b) Observations + assigned zeros

Figure 3: Density of points in the database. Color corresponds to a) number of observations in cell, or b) number of observations plus the number of assigned zero values.

that the location remains equatorward of the convection-reversal boundary under most conditions. The 182 plot shows significant scatter of velocity values for all IMF values, though the average shows a roughly 183 linear trend of increasing westward velocity with increasing negative B_z magnitude. The average val-184 ues of the velocity are offset from the highest density of points indicated by the color contours, which 185 shows that the distributions are non-Gaussian. Frame 4b shows that v_{ns} also has significant scatter, 186 however it remains small for all values of B_z . It demonstrates a nearly linear trend of increase with in-187 creasing negative B_z magnitude, however the trend is small and the spread of velocities is significantly 188 larger than the trend. 4c shows v_{ew} vs the auroral index A_l which is an indicator of substorm activ-189 ity. Again, there is significant scatter in the velocity values for all values of A_l . There is also a clear 190 nonlinear dependence of the velocity on the index. For small values of A_l , the velocity magnitude in-191 creases rapidly with increasingly negative A_l , while at higher index values the velocity increase is small. 192 Similar behavior is illustrated for the dependence of v_{ew} on the A_u index (4d). 193

Figure 5 is the same format as Figure 4 but for the bin at 81° latitude and 1300 MLT, which lies in the post-noon polar cap under most conditions. The dependencies differ significantly from the auroral zone dusk cell. v_{ns} at this location shows a clear nearly linear dependence on B_z , with positive (antisunward) values for negative B_z and negative (sunward) values for positive average B_z in excess of about 2.5 nT.



Figure 4: Dependence of velocity components in the bin at 67° magnetic latitude and 1800 MLT versus select parameters from the database. (a) shows the relationship between v_{ew} and the IMF z-component, (b) shows v_{ns} vs IMF z-component, (c) shows v_{ew} vs A_l , and (d) shows v_{ew} vs A_u .

The east-west velocity component is small magnitude and appears weakly correlated with the parameters $(B_z, A_u, \text{ and } A_l)$.

As the two figures show, the relationship between the velocity components and the database features is complex and varies from place to place. In some regions the velocity may be much more strongly correlated with one parameter than with another, while in another location the opposite is true. Because of this, none of the parameters was eliminated from consideration. All seven parameters were used to train the models.

The database was used to train an independent model of each velocity component $(v_{ns} \text{ and } v_{ew})$ in 205 each latitude-MLT grid cell. With the 17 latitude bins and 24 longitude bins, there are 408 grid cells. 206 Fitting the model components separately means that there are a total of 816 independent models. Three 207 algorithms were tested for forming the model. The algorithms were the LinearRegression, DecisionTreeRe-208 gressor, RandomForestRegressor provided by the Scikit-Learn software package (Pedregosa et al., 2011). 209 To test each algorithm, the data base was processed with each model in a ten-cell subset of the grid 210 space. To limit over-fitting, the maximum-depth hyperparameter for the Random Forest and Decision 211 Tree model was set to 15. The resulting models were used to predict velocities in a sample of data out-212 side of the training set and the model with the lowest root-mean-squared error (rmse) was selected. 213 In each case, the Random Forest Regressor was substantially better than the others. For example in 214 the bin at 67° latitude 1800 MLT, linear prediction of v_{ew} resulted in a rmse of 169.3 m/s, the decision 215 tree resulted in a rmse of 126.5 m/s, and the random forests resulted in a rmse of 113.2 m/s. In Scik-216 itLearn, the Random Forest model is trained by fitting multiple decision trees to random subsamples 217 of the inuput data and aggregating the predictions of all the trees. This is one way to address the over-218 fitting in addition to controlling the maximum depth of each tree. 219

After model selection on the subset of grid cells, the full dataset was split into data from the years 2014 220 to 2017, which was used to train the model, with data from 2013 used as a test set. Figure 6 shows 221 a 2000 sample interval from the model in the bin latitude 67°, MLT 1800. The horizontal axis is sam-222 ple number from the database which corresponds to time, however because of the data are not con-223 tinuous multiple time intervals contribute to the plot resulting in discontinuities in the plot traces that 224 do not represent temporal discontinuities of the values. The upper frame of the plot shows the observed 225 values in red, the model predictions in green, and the predictions from the TS18 model in blue. The 226 lower two frames show the IMF y and z components, and the A_{u} and A_{l} indicies. For most of the in-227

-11-



Figure 5: Same as Figure 4 except for the bin at 81° magnetic latitude and 1300 MLT.



Figure 6: Sampling of the driving features and model predictions for the bin at 67° magnetic latitude and 1800 MLT. The 2000 sample interval is composed of multiple time intervals. In the top panel, the red trace is the observed velocity, the green trace is the velocity predicted by the ML model, and the blue trace is the velocity predicted by the TS18 model.



Figure 7: Scatter of model predictions versus observed values of v_{ew} in the bin at 67° latitude 1800 MLT.

terval, the predicted value Is close to the observed value, with differences of less than 100 m/s. The predictions from the TS18 model at times show much larger differences from the observations as is especially well illustrated in the values from samples between 12400 and 13000 where the difference is on the order of 500 m/s.

Figure 7 shows the difference between the predictions and observations for the full year of 2013. The 232 figure shows the scatter of predicted velocity (vertical axis) versus measured velocity (horizontal axis) 233 in the bin at 67° and 1800 MLT. The solid black circles show the average values and the horizontal lines 234 show the average plus and minus one standard deviation. The average values follow the equality line 235 for most of the range, with significant deviation only for values where there are relatively few points. 236 Where the average values lie above the equality line, there is a small bias (< 50 m/s) for the model 237 values to be smaller magnitude than the observed values. While the scatter appears large, the stan-238 dard deviations demonstrate that the majority of the predictions are within 100 m/s of the observations 239 for all values with a significant number of observations. 240

-14-

Figure 8 shows the result of running the model for two intervals with similar IMF values but significantly 241 different values for A_u and A_l . The north-south at east-west models were run for each grid cell and 242 then combined to form vectors. Each grid cell is independent, so there was no guarantee that the out-243 put would produce a coherent convection pattern. The results do in fact illustrate a well defined co-244 herent convection pattern that is consistent with expectations based upon the observed driving con-245 ditions. During the two intervals the IMF was southward with $B_z = 4.7 \, \mathrm{nT}$ and $B_y = -3.85 \, \mathrm{nT}$ in the 246 first interval (8a) and $B_y = -0.76 \,\mathrm{nT}$ in the second interval (8b). In frame 8a, the auroral indicies are 247 small with $A_u = 84 \text{ nT}$ and $A_l = -31 \text{ nT}$. In frame 8b A_u was roughly two and a half times and A_l was 248 roughly four times the value in 8a . While the IMF values are similar, the patterns show significant dif-249 ferences. The most obvious of which are that the pattern driven by the larger values of A_u and A_l ex-250 tends to lower latitudes and has larger magnitude at nearly all locations. In 8a the convection is con-251 fined to latitudes above about 65°, while in 8b it extends to below 60° in the pre-midnight sector. Day-252 side plasma flows extend to slightly lower latitude in the later interval, though not by as much as the 253 night-side flows. The main difference in the dayside is that the direction and local time of plasma en-254 try to the polar cap reflects the influence of the larger IMF B_y in the earlier interval. The nightside exit 255 of plasma differs significantly between the two plots. In 8a flow near midnight is small magnitude and 256 mainly equatorward before turning to connect to the return flow regions. In 8b the dawn cell is roughly 257 "D" shaped with the flow turning directly from cross-cap to the return flow region, while in the dusk 258 cell, there is the flow rotates first dawnward before rotating back to connect with the dusk return flow. 259 This dusk-cell shear flow illustrates the development of the Harang Discontinuity with increasing au-260 roral activity. 261



Figure 8: Output of the 816 independent models displayed on the latitude-MLT grid for similar IMF conditions but differing auroral indicies a) $B_z = -4.53 \text{ nT}$, $B_y = -2.45 \text{ nT}$, $A_u = 82 \text{ nT}$, $A_l = -34 \text{ nT}$ b) $B_z = -4.6 \text{ nT}$, $B_y = -1.41 \text{ nT}$, $A_u = 157 \text{ nT}$, $A_l = -128 \text{ nT}$

To examine the accuracy of the model prediction over the entire grid, the root-mean-squared difference between the predictions and observations were calculated in each grid cell for all observations the year of 2013. The models were used to predict the velocity components at each time for which there as an observation in a given grid cell based upon the values in the OMNI database from the time of the observation. For comparison, the TS18 model was used similarly to predict the velocity components and compared to the observations. Figure 9 shows the results for v_{ew} over the grid for both the ML model and the TS18 model.

manuscript submitted to Space Weather



Figure 9: Root-mean-squared difference between model predictions and the observations of the eastwest velocity component for the year of 2013. a) RMSE from the ML model, b) TS18 model

The figure shows that the RMSE for the ML model is less than about 250 m/s over the entire grid. The lowest values occur on the low-latitude dayside, where velocities are typically low. The highest values occur in the prenoon sector between 75° and 80°, and at the lowest latitude bins near dawn and dusk. In addition errors on the night-side are highest in the region between 70° and 75°. The plot for the TS18 model shows higher RMSE for all bins, with particularly large errors (>350 m/s) near 70° for all local times.

275 **3 Discussion**

The need for accurate forecasting of space weather increases on a nearly daily basis. There isn't a better example of this than the requirement for accurate orbit prediction that becomes more critical with the launch of every new low-Earth-orbit satellite. Orbit prediction is based on thermospheric density, which can be predicted using global circulation models driven by convection models such as described in this study. Hence, it is imperative that we have models of convection that accurately capture the variation of the high-latitude potential with IMF and internal magnetospheric state. The ML model pre-



Figure 10: Potential pattern resulting from a spherical harmonic expansion constrained by the ML model

sented here shows a marked improvement over traditional climatological models, however it requires specification of auroral indicies, A_l and A_u , which are based upon magnetometer observations. Recently there have been successful efforts at predicting these and other magnetospheric indicies using machine learning techniques using the solar wind and IMF as inputs. With these predicted indicies it would be possible to predict the convection, and in turn the thermospheric density, several hours into the future (Topliff et al., 2019).

For retrospective studies, observations of the indicies are readily available and can be used to select model 288 patterns for use in GCMs or to serve as a constraint on instantaneous convection patterns generated 289 using MapPotential. GCMs and MapPotential use models of the electrostatic potential rather than ve-290 locities. Such maps can be generated from the output of the ML model by using the same technique 291 that has been applied in generating other climatological models. Figure 10 shows the result of expand-292 ing the potential in spherical harmonics using the ML model velocity field in Figure 8b as a constraint 293 on the fit. Because the fit is a functional expansion it can be calculated on a fine grid, which gives the 294 smooth variation with position illustrated in the figure. The pattern is similar to that shown in Figure 295 1b, which was generated from TS18 using similar values of the IMF. The dusk cell is not quite as round 296 in the ML model plot and the flow near midnight extends to lower latitudes. 297

-18-

Figure 8, shows that changes in A_l and A_u result in changes in the convection predicted by the ML 298 model which reflect the expected behavior of the polar cap. Convection extended to lower latitudes in 299 Figure 8b than in Figure 8a in response to the significantly larger values of A_u and A_l . The latitude 300 of the convection reversal is impacted by the changes in the indicies, though the change is not uniform 301 in local time. At dawn the boundary was lower by several degrees in 8b than in 8a, and the develop-302 ment of the Harang Discontinuity appears in 8b, which gives a convection reversal at low latitudes in 303 the premidinght region and extending across 0000 MLT. The extension of the convection reversal across 304 midnight results in the tongue of negative potential extending across midnight illustrated in Figure 10. 305 The dayside convection reversal and that at dusk shows little difference between the two intervals. The 306 dayside differences between the two intervals that do appear are more likely due to the differing val-307 ues of the IMF B_u . 308

As illustrated by the comparisons to the TS18 climatology, the ML model represents an improvement 309 over models that do not attempt to capture variability driven by internal magnetospheric processes. Fig-310 ure 9 showed that the RMSE of the ML model predictions vs SuperDARN observations is substantially 311 lower than that for the TS18 model. The region of large RMSE in the TS18 model concentrated be-312 tween 70° and 75° is most likely due to the inability of the model to capture the expansion and con-313 traction of the polar cap boundary during substorm cycles. The region of large RMSE is close to the 314 average position of the polar cap boundary, which is close to the latitude of the convection reversal bound-315 ary (CRB). The CRB is of course, the latitude separating the antisunward flow in the polar cap and the 316 sunward flow on field lines that map into the magnetosphere. When the polar cap expands so that the 317 CRB latitude is below it's average position for a given set of IMF/solar wind conditions, grid cells just 318 below the boundary would be predicted to lie in magnetosphere and have sunward flow, while in fact 319 they lie in the polar cap and have antisunward flow. Likewise, when the boundary contracts above it's 320 average position, grid cells just poleward of the average boundary would be predicted to have antisun-321 ward flow while in fact it is sunward. When either of these conditions happens, the difference between 322 the prediction and the observation would be on the order of double the average magnitude of the ve-323 locity in the cell. Since the boundary is so dynamic it is likely that this is a common occurrence, which 324 is reflected by the large average errors that appear in the figure. It should be noted that the largest er-325 rors in the ML model also appear in this region, however they are significantly lower magnitude than 326 in the TS18 model indicating that the ML model does a better job of representing the changes in the 327 boundary position. 328

-19-

Another advantage of the ML model over existing climatologies is that the way it was formed allows 329 for easy characterization of the distribution of velocities in each grid cell, which can be used when as-330 similating the model output. The RMSE is returned for each grid cell as part of the ML regression. If 331 the model is assumed to be unbiased, the RMSE value can be used as the square root of the variance 332 and combined with the model output value to generate a distribution function assuming a Gaussian dis-333 tribution. A Bayesian assimilation scheme would use the distribution of model as prior information. Hav-334 ing an RMSE in each grid cell contrasts with the information obtained when expanding the potential 335 in orthogonal functions. In such a fit, the function coefficients are returned, which distributes errors 336 over the domain. While it would be possible to generate a covariance matrix for a functional expan-337 sion, it requires the extra step of using the model to predict a local velocity and calculating the sam-338 ple variance around that value. 339

340 **4** Conclusions

This paper describes a climatological model of high-latitude convection derived using machine learn-341 ing (ML) techniques applied to observations from the SuperDARN radar network. The model was gen-342 erated from a database of four years of observations and tested over a separate fifth year. SuperDARN 343 convection patterns were generated for every five minute period over the five year period. From those 344 patterns, velocity vectors were calculated at locations where there was at least one radar contributing 345 observations. Those velocities were separated into north-south, and east-west components and sorted 346 into a magnetic local time - magnetic latitude grid that ran from 55° to the magnetic pole with a bin 347 size of 2°, and MLT bins of 1-hour. 348

In each MLT-MLAT bin, the two velocity components were used separately to train a ML model using random forests regression. Random forests was selected after testing three different ML algorithms to find the one that produced the lowest RMSE in a subset of the points in the grid. The features used to train the model were the IMF components B_x , B_y , and B_z ; the solar wind velocity, v_{sw} ; the auroral indicies, A_u and A_l ; and the geomagnetic index, *SYM-H*.

After the model was trained on data from the years 2014 to 2018 (inclusive), it was tested using data from the year 2013. Predictions from the model were compared to the SuperDARN observations and distributions of predicted versus observed velocity were examined. While there was significant scatter

-20-

of the predictions around the line of equality with the observations, the average of the distributions tracked

 $_{\scriptscriptstyle 358}$ the average measured velocity well with a small bias to lower values. The standard deviation of the model

³⁵⁹ predictions was less than 100 m/s for all bins where there were a significant number of observations.

RMSE values for the model were compared to those from the TS18 model in each bin of the grid. The ML model exhibited smaller errors than TS18 at all locations. In particular, errors in the ML showed the largest improvement over TS18 in bins that are near the average latitude of the convection reversal boundary. It is likely that the improvement was due to the ML model's ability to expand and contract in latitude in response to changes of A_l and A_{ll} .

The software for generating the model is free and available for download from the scikit-learn web site. The web site has links to numerous examples and tutorials for application of the various algorithms it provides. The software is simple to use even by senior investigators with no prior experience with ML techniques. Despite the simplicity, good results can be obtained with some time spent reading the tutorials.

370 **5 Data Availability**

The raw SuperDARN data are available from the British Antarctic Survey (BAS) SuperDARN data server (https://www.bas.ac.uk/project/superdarn).

Acknowledgments: This work is supported by the Defense Advanced Research Projects Agency (DARPA) 373 through US Department of the Interior award D19AC00009 to the Georgia Institute of Technology and 374 subaward to The Pennsylvania State University. SuperDARN operations and research at Pennsylvania 375 State University are supported under NSF Grants PLR-1443504 from the Office of Polar Programs, and 376 AGS-1934419 from the Geospace Section of NSF Division of Atmospheric and Geospace Sciences. The 377 authors acknowledge the use of SuperDARN data. SuperDARN is a collection of radars funded by na-378 tional scientific funding agencies of Australia, Canada, China, France, Italy, Japan, Norway, South Africa, 379 United Kingdom and the United States of America. We acknowledge use of NASA/GSFC's Space Physics 380 Data Facility's OMNIWeb service, and OMNI data. 381

manuscript submitted to Space Weather

382 References

383	Bristow, W. A., & Jensen, P. (2007). A superposed epoch study of superdarn convection observa-
384	tions during substorms. Journal of Geophysical Research: Space Physics, 112(A6). Retrieved
385	<pre>from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2006JA012049</pre>
386	doi: 10.1029/2006JA012049
387	Foster, J. C. (1983). An empirical electric field model derived from chatanika radar data.
388	Journal of Geophysical Research: Space Physics, 88(A2), 981-987. Retrieved from
389	https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JA088iA02p00981
390	doi: 10.1029/JA088iA02p00981
391	Greenwald, R. A., et al. (1995). (1995), DARN/SuperDARN: A global view of high-latitude con-
392	vection. Space Sci. Rev, 71, 763-796.
393	Heelis, R. A., Lowell, J. K., & Spiro, R. W. (1982). A model of the high-latitude ionospheric
394	convection pattern. Journal of Geophysical Research: Space Physics, 87(A8), 6339-6345.
395	Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/
396	JA087iA08p06339 doi: 10.1029/JA087iA08p06339
397	Heppner, J. P., & Maynard, N. C. (1987). Empirical high-latitude electric field models.
398	Journal of Geophysical Research: Space Physics, 92(A5), 4467-4489. Retrieved from
399	https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JA092iA05p04467
400	doi: 10.1029/JA092iA05p04467
401	King, J. H., & Papitashvili, N. E. (2005). Solar wind spatial scales in and comparisons of hourly
402	wind and ace plasma and magnetic field data. Journal of Geophysical Research: Space
403	<i>Physics</i> , 110(A2). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/
404	abs/10.1029/2004JA010649 doi: 10.1029/2004JA010649
405	Papitashvili, V. O., Belov, B. A., Faermark, D. S., Feldstein, Y. I., Golyshev, S. A., Gro-
406	mova, L. I., & Levitin, A. E. (1994). Electric potential patterns in the northern
407	and southern polar regions parameterized by the interplanetary magnetic field. Jour-
408	nal of Geophysical Research: Space Physics, 99(A7), 13251-13262. Retrieved from
409	https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/94JA00822 doi:
410	10.1029/94JA00822
411	Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Duchesnay, E.
412	(2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12,
413	2825–2830.

-22-

- Ridley, A. J., Deng, Y., & Tóth, G. (2006, May). The global ionosphere thermosphere model. 414 Journal of Atmospheric and Solar-Terrestrial Physics, 68(8), 839-864. doi: 10.1016/j.jastp 415 .2006.01.008 416
- Roble, R. G., & Ridley, E. C. (1994). A thermosphere-ionosphere-mesosphere-electrodynamics 417 general circulation model (time-gcm): Equinox solar cycle minimum simulations (30–500 418 Geophysical Research Letters, 21(6), 417-420. Retrieved from https://agupubs km).
- .onlinelibrary.wiley.com/doi/abs/10.1029/93GL03391 doi: 10.1029/93GL03391 420

419

- Ruohoniemi, J. M., & Baker, K. B. (1998). Large-scale imaging of high-latitude convection with 421 Super Dual Auroral Radar Network HF radar observations. J. Geophys. Res, 103, 20,797. 422
- Ruohoniemi, J. M., & Greenwald, R. A. (1996). Statistical patterns of high-latitude convection 423 obtained from goose bay hf radar observations. Journal of Geophysical Research: Space 424 Physics, 101(A10), 21743-21763. Retrieved from https://agupubs.onlinelibrary 425 .wiley.com/doi/abs/10.1029/96JA01584 doi: 10.1029/96JA01584 426
- Schulz, M. (1997). Direct influence of ring current on auroral oval diameter. Journal of Geo-427 physical Research: Space Physics, 102(A7), 14149-14154. Retrieved from https://agupubs 428 .onlinelibrary.wiley.com/doi/abs/10.1029/97JA00827 doi: 10.1029/97JA00827 429
- Shepherd, S. G., & Ruohoniemi, J. M. (2000).Electrostatic potential patterns in the high-430 latitude ionosphere constrained by superdarn measurements. Journal of Geophys-431 ical Research: Space Physics, 105(A10), 23005-23014. Retrieved from https:// 432 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2000JA000171 doi: 433 https://doi.org/10.1029/2000JA000171 434
- Siscoe, G. L., & Huang, T. S. (1985).Polar cap inflation and deflation. Journal of 435 Geophysical Research: Space Physics, 90(A1), 543-547. Retrieved from https:// 436 agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/JA090iA01p00543 doi: 437 10.1029/JA090iA01p00543 438
- Thomas, E. G., & Shepherd, S. G. (2018). Statistical patterns of ionospheric convection 439 derived from mid-latitude, high-latitude, and polar superdarn hf radar observations. 440
- Journal of Geophysical Research: Space Physics, 123(4), 3196-3216. Retrieved from 441 https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2018JA025280 doi: 442 10.1002/2018JA025280 443
- Topliff, C., Cohen, M., & Bristow, W. (2019). Simultaneously forecasting global geomagnetic ac-444 tivity using recurrent networks. Machine Learning and the Physical Sciences Workshop, Ad-445 vances in Neural Information Processing Systems. 446

-23-

- Weimer, D. R. (2005). Improved ionospheric electrodynamic models and application to calculating
 joule heating rates. Journal of Geophysical Research: Space Physics, 110(A5). Retrieved
 from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2004JA010884
- 450 doi: 10.1029/2004JA010884