

Spatiotemporal Relational Random Forests

Timothy A. Supinie
School of Meteorology
University of Oklahoma
Norman, OK 73072
Timothy.A.Supinie-1@ou.edu

Amy McGovern
School of Computer Science
University of Oklahoma
Norman, OK 73019
amcgovern@ou.edu

John Williams, Jennifer Abernethy
Research Applications Laboratory
National Center for Atmospheric Research
Boulder, CO 80301
jkwillia@ucar.edu, aberneth@ucar.edu

Abstract—We introduce and validate **Spatiotemporal Relational Random Forests**, which are random forests created with spatiotemporal relational probability trees. We build on the documented success of random forests by bringing spatiotemporal capabilities to the trees, enabling them to identify critical spatial, temporal, and spatiotemporal features in the data. We validate our results on simulated data and real-world convectively-induced turbulence data from a commercial airline flying in the continental United States.

Keywords-Spatiotemporal data mining, Relational learning, Random forests, Turbulence

I. INTRODUCTION

The contribution of this paper is the introduction of spatiotemporal relational random forests (SRRFs). These are random forests that directly handle spatiotemporal relational data. Random forests as introduced by Breiman [1] use the standard C4.5 decision trees [2] to create the individual trees, with random sampling of data and splitting attributes. SRRFs use spatiotemporal relational probability trees [3], which are relational probability trees [4] that handle spatially and temporally varying relational data.

This work is motivated by two real-world domains where space and time are critical to effective understanding and classification. The first domain is predicting convectively induced turbulence (CIT) for aviation in North America. CIT refers to turbulence generated from thunderstorms as opposed to clear air turbulence or mountain wave turbulence. Empirical data and modelling studies suggest that thunderstorms generate nearby turbulence on a short temporal scale. Gravity waves and outflow winds may also produce turbulence much farther away over a longer time scale [5], [6], [7]. The second motivating domain is tornadoes and severe thunderstorms. In both domains, attributes of the weather vary over space and time and these spatiotemporal relationships are critical for understanding and predicting the evolution of the domain. In this paper, we apply SRRFs to the turbulence domain but in separate work by some of the authors [3], we are applying both SRPTs and SRRFs to simulated tornadic storms.

Meteorologists study such data using high-level features and the relationships between these features. For example, in predicting CIT, the distance to the nearest thunderstorm

and the intensity and relative altitude of that storm are critical, along with other environmental features [6], [8]. By focusing on the relational nature of the data, we maintain the ability to reason both about attributes and about the critical spatiotemporal relationships among the features themselves.

Random forests have been successful in a variety of applications (e.g., [9], [10], [11], [12]) and are well analyzed and understood (e.g., [13], [14]). Prior work in the turbulence domain demonstrated that random forests were a promising approach [8], [15]. The SRRF approach differs from many other methods in spatiotemporal analysis and pattern discovery. One of the biggest differences is that we use relational data, which enables us to find spatiotemporal patterns at the object and relationship level. For example, we can identify patterns involving the temporal changes in spatial relationships such as two objects moving past each other. Temporal relational data mining is a new area [3], [16], [17] and the spatial component is relatively unexplored. Another difference between our approach and other non-relational spatiotemporal data mining methods such as [18], [19], [20], [21], [22] is the focus on human-readable models.

II. APPROACH

Standard random forests are constructed by creating k C4.5 trees, each of which is trained on a different subset of the data and with a random subset of attributes at each node. Each subset of the data is created by sampling a set of size n with replacement from the original training set of size n . For Spatiotemporal Relational Random Forests, we use Spatiotemporal Relational Probability Trees (SRPTs) [3]. SRPTs are probability estimation trees designed for spatiotemporal relational data. They are an extension of the Relational Probability Trees [4], which were designed for static relational data. Probability estimation trees are decision trees with probabilities at the leaf nodes.

The training data for SRPTs and SRRFs are represented as a spatiotemporal attributed graph. The schema for our data is shown in Figure 1. Objects are represented by vertices in the graph and relations by edges. Both objects and relations can have attributes associated with them. The object types in our turbulence data are: aircraft, regions of vertically integrated liquid (VIL), and regions of “rain”,

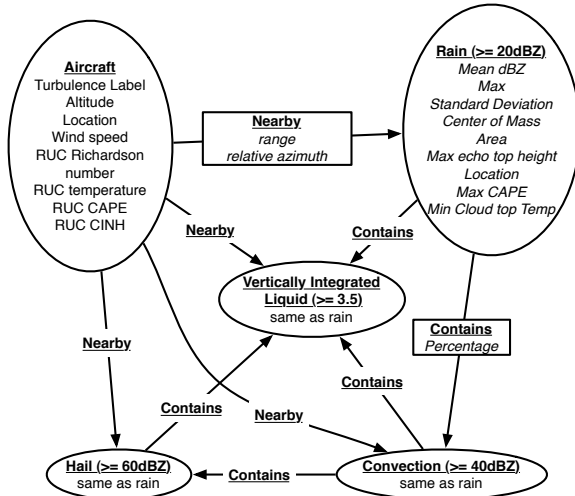


Figure 1. Schema for the turbulence data. Objects are represented with ellipses and relations with lines and rectangles. Temporal attributes are italicized.

“convection” (thunderstorms), and “hail” derived from a composite radar reflectivity mosaic. Each of these can be “nearby” or “containing” one another. The attributes on all objects and relations except for the aircraft are temporal. The aircraft has static attributes and the goal is to predict the “turbulence label” attribute.

SRPTs are built using the standard greedy decision tree approach with the chi-squared statistical test used to score new nodes. As with the trees for standard random forests, SRPTs are not pruned. At each decision node, an SRPT asks a question about the spatiotemporal relational graph. The questions cover a variety of spatial, temporal, and spatiotemporal changes in the graph. Example questions include “does a convective region exist nearby the aircraft?” and “is there a hail region with area $\geq 10 \text{ km}^2$ within 10 km of the aircraft?” The full set of question types can be found in [3].

Because the search space for the set of possible questions at each node is so large, the SRPTs are constructed by stochastically sampling the set of all possible distinctions. The number of samples is controlled through a user-specified parameter. Although higher numbers of samples yield better trees, the computational effort extracts a price in the training time of the trees. This sampling parameter is similar to the parameter in random forests that controls the number of attributes examined by C4.5 for each node in a tree and we treat it as such. We empirically examine the effect of this parameter.

III. CONVECTIVELY INDUCED TURBULENCE

Pilots’ ability to avoid turbulence during flight affects the comfort and safety of millions of airline passengers annually. Of all weather-related commercial aircraft incidents, 65% can be attributed to turbulence encounters, and major carriers

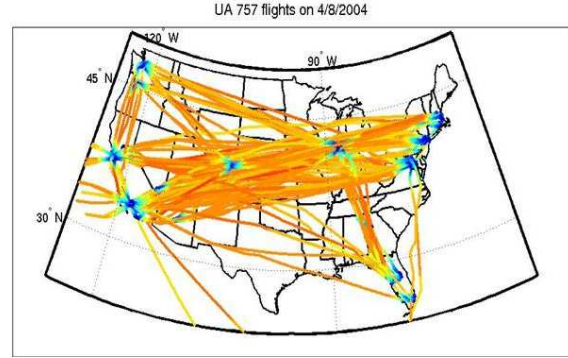


Figure 2. The geographical distribution of automated EDR measurements from United Airlines 757 aircraft over a 24-hour period.

estimate that they receive hundreds of injury claims and pay out “tens of millions” of dollars per year [23]. Among the different types of turbulence, turbulence in and around thunderstorms, which is termed convectively-induced turbulence (CIT), may be responsible for over 60% of turbulence-related aircraft accidents [24].

Mechanisms for the generation and propagation of CIT are not currently well-understood by researchers. CIT is commonly thought to be related to the proximity (vertical and horizontal) to convection and to its intensity, depth and extent, but is also known to propagate through gravity waves that may break above or away from the storm [25]. This may also be produced by anvil winds interacting with the surrounding environment [7]. Due to these uncertainties, Federal Aviation Administration (FAA) guidelines require aircraft to circumnavigate thunderstorms by wide margins both horizontally and vertically to mitigate the risk of encountering dangerous turbulence, making large regions of airspace unavailable to aircraft on days of widespread convection.

An active area of research and development is the enhancement of the FAA’s clear-air turbulence forecast (Graphical Turbulence Guidance; [23]) with in-cloud turbulence detection and near-cloud turbulence diagnosis capabilities [26], [15]. As part of this effort, random forests have been used to evaluate the potential predictive value of numerical weather prediction (NWP) model and observation data and to develop an empirical model for producing high-resolution, rapid-update, 3D probabilistic assessments of light, moderate, and severe turbulence. This approach relies heavily on using automated *in situ* turbulence measurements from commercial aircraft to develop and calibrate the model. Despite some success using this approach, it is expected that an improved methodology for the analysis and exploitation of spatiotemporal relationships between turbulence and storm characteristics and environmental features could provide improved CIT prediction.

The high temporal and spatial resolution of *in situ* aircraft

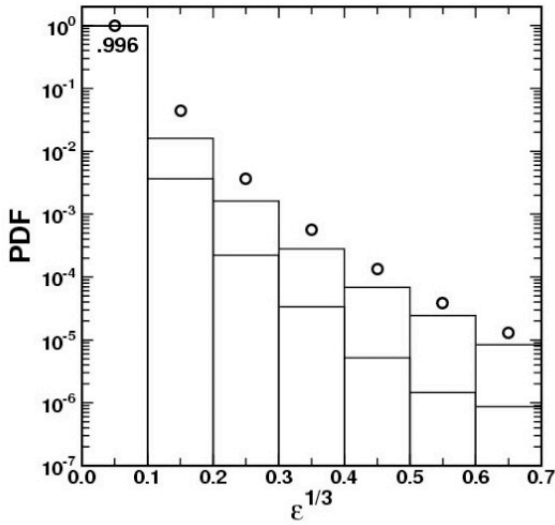


Figure 3. Distribution of binned 1-minute median (lower bar) and peak, i.e. 95th percentile, (upper bar) EDR ($\epsilon^{1/3}$) values measured by United Airlines Boeing 757 aircraft over a three month time period. The open circles are estimates based on an assumed lognormal distribution with parameters derived from an operational NWP model [27]. The difference may reflect the ability of commercial air carriers to successfully avoid turbulence. Taken from [23].

measurements of eddy dissipation rate (EDR, a measure of turbulence), made every minute by some commercial aircraft, make them good “truth” data for developing and evaluating CIT prediction algorithms. Figure 2 shows the geographical distribution of EDR reports from United Airlines over a 24-hour period. Elevated turbulence is a rare phenomenon, however [23], and aircraft reports do not provide representative sampling of the atmosphere. Figure 3 shows the distribution of EDR measurements and a theoretical distribution based on an analysis of numerical weather prediction model output. The measured frequency of elevated turbulence is lower than the theoretical estimate, presumably because aircraft are able to use available forecasts and pilot report information to avoid some turbulence. Thus, these turbulence measurements likely represent primarily unexpected turbulence encounters, making the prediction problem even more challenging. The United EDR reports consist of both a peak EDR and median EDR measured over a one-minute time period, binned into one of seven intensity bins that may be roughly characterized as representing null, light, light-to-moderate, moderate, moderate-to-severe, severe, and extreme turbulence, respectively. For this initial study, we consider the binary classification problem of predicting whether turbulence is moderate-or-greater (bin 4 or higher). We use EDR measurements from a three month period over the summer of 2007.

IV. DATA SETS

The **BlobWorld** is a simulated data set with spatiotemporal characteristics that we use as a test case. The world consists of agents of two types randomly distributed on a grid and a central agent whose type is unknown. The grid and the number of agents are parameterized. For these results, we used a 5 by 5 grid with 5 agents. All agents are assigned a “type” attribute when they are created and all agents have a “happiness” attribute, the value of which is dependent on the inverse square of its distance from agents of the same type. The central agent moves pseudo-randomly on the grid, with a bias toward agents of its type, and updates its happiness at each time step. The task of the random forest is to determine the type of the central agent given the spatiotemporal behavior of this agent with respect to the other agents of known type. This is a difficult task but particularly well suited to spatiotemporal models.

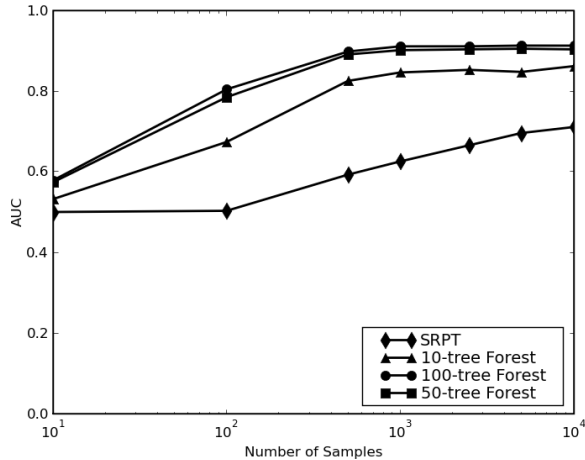
The **turbulence** data described above may be resampled in multiple ways for use in the SRRFs: by distribution of turbulence intensities, geographically, and by proximity to convection. The overwhelming number of “no turbulence” data points in the full data set would cause the algorithm to achieve a 99% accuracy rate by exclusively predicting “no turbulence.” In order to make a valid predictor for a rare phenomenon, we rebalanced the data set to retain roughly 3% of the “no turbulence” EDR data points. To eliminate other sources of turbulence, we used only data points within the continental U.S. (CONUS), above 15000 ft in altitude, and within 40 nautical miles of convection.

The graphs used a combination of meteorological variables collocated with the *in situ* aircraft measurements along with NWP model-based convective available potential energy (CAPE), convective inhibition (CIN), Richardson number, temperature, and composite radar reflectivity data and radar-derived vertically-integrated liquid (VIL). The “aircraft” object was given the attributes as shown in Figure 1. The “rain,” “convection,” “hail,” and “VIL” objects were constructed from the composite reflectivity data by the application of a simple threshold. “Rain” objects are connected areas of radar reflectivity above 20 dBZ, likewise “convection” objects above 40 dBZ, “hail” objects above 60 dBZ, and “VIL” above 3.5. The threshold values were chosen because they are simple rules of thumb for categorizing radar echoes. After applying the thresholds and retrieving connected areas within 40 nautical miles, the attributes given in Figure 1 were computed.

V. PRELIMINARY RESULTS

A. BlobWorld

We use BlobWorld to examine the effect of the main forest parameters on the performance of the forest. We varied the number of trees in the forest and the distinction sampling rate. Forests of 10, 50, and 100 trees were used along



(a) Performance of the forest

Figure 4. Average AUC for BlobWorld as a function of sample size and number of trees in the forest.

with sampling rates of 10, 100, 500, 1000, 2500, 5000, and 10000 samples per tree node. For comparison, we also examined the performance of an individual SRPT, varying only the number of samples. Each combination of these parameters was run 30 times and the area under the Receiver Operating Characteristic (ROC) curve (AUC) was computed and averaged over all the runs.

Figure 4 shows the results of varying the number of trees in the forest and the number of samples for each tree. As expected, performance of the forest and SRPT improves as the number of samples used to grow the tree increases. However, performance for all three forest sizes asymptotes as the number of samples reaches approximately 1000. This is to be expected, and similar behavior was observed for standard random forests. [1].

Figure 4 also shows that the performance increases as a function of the number of trees in the forest but that the gain asymptotes as the forest gets larger and larger. We used a single-tail paired t -test to compare the performance across sets of parameters and we corrected the p-values using a Bonferroni adjustment [28] to obtain an effective p-value of 0.01. We compared the performance of forests of 10 and 100 trees and of the SRPT against the 100 tree forest. The difference for both sets of comparisons is statistically significant at all levels of sampling.

Figure 4 visually shows that the two parameters (forest size and sampling) both have an effect on performance, which we verified statistically using an ANOVA test. As the table shows, not only do the two individual parameters show a statistically significant effect on the performance of the algorithm, but there is a significant interaction effect between the parameters. This is to be expected since the sampling size improves the quality of the trees and the number of

Factor	Blob	Turbulence
Forest Size	0	0
Number of Samples	0	0
Forest Size \times Samples	0	0.63

Table I

ANOVA P-VALUES FOR SIGNIFICANCE OF PARAMETER EFFECTS ON THE SIZE OF THE FOREST AND THE NUMBER OF SAMPLES USED TO BUILD THE TREES. STATISTICALLY SIGNIFICANT RESULTS ARE BOLDED.

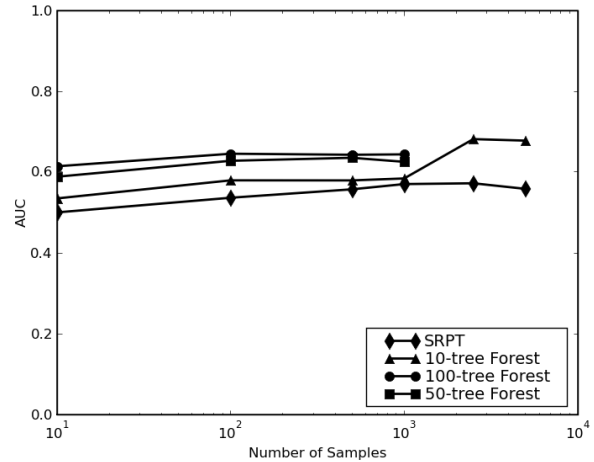


Figure 5. Average AUC for the subset of turbulence data over 30 runs as a function of sample size and number of trees in the forest.

trees improves the size of the ensemble. We repeated the ANOVA test for the performance of the individual SRPTs for the number of samples and it was also statistically significant ($p = 0$), which verifies the hypothesis that the performance of the SRPTs improves as a function of the number of samples.

B. Turbulence

Our initial turbulence experiments used only a few variables and relations, all centered on the aircraft. The results were not as promising as we would have liked so we examined the data and added additional meteorological variables, an improved graph structure, including “contains” relations and additional relationships between all objects (as opposed to just those that centered on the aircraft), and we resampled the data set to be more balanced. Because processing of this data had not been completed by the time of final paper submission, we ran the tests on 8 days of training data and 8 days of test data. This yielded 117 training graphs and 236 testing graphs. We varied the forest size from 10, 50, and 100 and varied the sampling rates from 10, 100, 500, and 1000. We also used sampling rates of 2500 and 5000 for the 10-tree forests and the SRPTs but could not do so for the 50 and 100 tree forests due to time constraints.

Figure 5 shows the average AUC across 30 runs of the SRRF on the turbulence data. The best performance is 0.68 with 10-tree forests at a sampling rate of 2500, which is

very promising. We theorize that the performance would be better with higher numbers of trees and we are working on these results currently. The unexpected jump in performance at the highest levels of sampling is most likely due to the difficulty of the problem. At lower levels of sampling, the individual trees are unable to identify high quality concepts.

As with the BlobWorld results, higher sample sizes and higher numbers of trees appear to have a positive effect on the performance of the forest. We verified this effect using the ANOVA test but, because ANOVA requires a full matrix, we could only verify this up to 1000 samples. The results of this test are given in Table I. We found that both the sample size and the size of the forests have a statistically significant effect, though in the turbulence world, there is no statistically significant interaction between the two parameters. This result was unexpected because it does not correlate with the behavior of the BlobWorld results. We theorize it is because of the low scores on the turbulence data up to sampling sizes of 1000 and we expect that the addition of the 50 and 100 tree forests at higher sampling values will change this result.

We also compared the results of the forests with the performance of a single SRPT. This comparison is also shown in Figure 5. Using the Bonferroni adjusted p-value described above and a single-tailed paired t -test, the performance of the SRPT and the SRRF are statistically different for all levels of sampling for the 10, 50 and 100 tree forests. This performance includes the higher sampling values for the 10-tree forest and the SRPT. We also compared the performance of the 10 tree and 100 tree forests, repeating the same single-tailed paired t -test with the Bonferroni adjusted p-value. The test indicates that they are statistically significant for all levels of sampling, meaning the performance of the 100 tree forest is statistically better than the performance of the 10 tree forest.

Although the initial performance of the SRRFs on the turbulence data was not as strong as on the BlobWorld data, we believe the primary issues are lack of training data and too few samples to produce a tree of reasonable quality. Due to the size of the data, it takes a significant amount of time to pre-process and we were not able to process the entire three months of turbulence data in time for the final publication of the paper. We expect that the results will improve significantly with this final data and with additional sampling and we are continuing this work.

VI. CONCLUSIONS

We have introduced Spatiotemporal Relational Random Forests, a novel approach to random forests that enable the forests to work with complex spatial and temporal relational data sets. We demonstrated that the random forests are a viable approach on the simulated BlobWorld data. We focused on only a single real-world application domain but we believe the applicability of this approach will be broad

and we are exploring a wide variety of additional real-world applications in related work.

Our goal is to refine the SRRF model and incorporate it into a CIT prediction algorithm designed to enhance the FAA's Graphical Turbulence Guidance capability. An important consideration will be making the algorithm efficient, which will require additional parameter exploration and identification of the best parameters for predicting CIT. Although we hand-picked meteorological variables for inclusion in the SRRF in the current results, random forests can be used to analyze the importance of various attributes and variables. We are in the process of implementing this importance analysis and expect that it will improve our ability to simplify the model by reducing the number of required input variables, which will improve SRRF performance and increase the applicability of the method.

One problem with decision tree approaches applied to real-world data sets is the known brittleness of the tree-based solutions [29]. Brittleness is defined as changes in the tree with little to no parameter variation in the learning algorithm. When applying these approaches to real-world domains, domain scientists often are wary to trust the answers of a such an approach. By moving to random forests, the overall stability of the algorithm is improved as well as the quality of the predictions. Although the ability to analyze individual trees is diminished as the size of the forest grows, approaches such as the importance analysis increase the utility for domain scientists. In turn, this will increase the applicability of the approach.

Another application for our future work is to expand the prediction of turbulence from only CIT to other forms of turbulence, such as clear air turbulence, which is expected to be related to proximity to the jet stream and other upper-level fronts and atmospheric features. Another path we are exploring is to apply the SRRF to the thunderstorm domain in order to better predict tornadoes. SRPTs have been applied to a data set of several hundred simulated thunderstorms with the goal of understanding the cause of the more severe storms [3]. This work was promising but the brittleness of the trees has proven problematic and we believe the random forests will address these issues.

ACKNOWLEDGEMENTS

This material is based upon work supported by the National Science Foundation under Grant Nos. IIS/CAREER/0746816, IIS/REU Supplement/0840956 and 0938138. This research was also supported by NASA under Grant No. NNS06AA61A. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Aeronautics and Space Administration or the National Science Foundation. We are grateful to Gary Blackburn, Jason Craig, Greg Meymaris and Bob Sharman

of NCAR for their contributions to compiling the turbulence dataset used in this study.

REFERENCES

- [1] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [2] J. R. Quinlan, *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
- [3] A. McGovern, N. Hiers, M. Collier, D. J. Gagne II, and R. A. Brown, "Spatiotemporal relational probability trees," in *Proceedings of the 2008 IEEE International Conference on Data Mining*, Pisa, Italy, 2008, pp. 935–940.
- [4] J. Neville, D. Jensen, L. Friedland, and M. Hay, "Learning relational probability trees," in *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2003, pp. 625–630.
- [5] R. Fovell, R. Sharman, and S. Trier, "A case study of convectively-induced clear air turbulence," *12. Conf. on Mesoscale Processes, Waterville Valley, NH (USA), 6-9 Aug 2007*, vol. Proceedings, p. paper 13.4, August 2007.
- [6] T. Lane and R. Sharman, "Some influences of background flow conditions on the generation of turbulence due to gravity wave breaking above deep convection," *J. Appl. Meteorol. Climatol.*, vol. 47, no. 11, pp. 2777–2796, November 2008.
- [7] S. Trier, R. Sharman, and R. Fovell, "Case studies of widespread turbulence in the vicinity of mesoscale convective systems using 13-km ruc analyses," *13. Conf. on Aviation, Range and Aerospace Meteorology (88th AMS Annual Meeting), New Orleans, 20-24 Jan 2008*, vol. Proceedings, p. paper 9.5, January 2008.
- [8] J. K. Williams, J. Craig, A. Cotter, and J. K. Wolff, "A hybrid machine learning and fuzzy logic approach to CIT diagnostic development," in *Preprints of the Fifth Conference on Artificial Intelligence Applications to Environmental Science*. American Meteorological Society, 2007, p. CDROM 1.2.
- [9] A. Bosch, A. Zisserman, and X. Munoz, "Image classification using random forests and ferns," in *Proceedings of the International Conference on Computer Vision*, 2007.
- [10] P. O. Fislason, J. A. Benediktsson, and J. Sveinsson, "Random forests for land cover classification," *Pattern Recognition Letters*, vol. 27, no. 4, pp. 294–300, 2006.
- [11] M. R. Segal, "Machine learning benchmarks and random forest regression," Center for Bioinformatics and Molecular Biostatistics, Tech. Rep., April 14 2004.
- [12] J. Williams, D. Ahijevych, S. Dettling, and M. Steiner, "Combining observations and model data for short-term storm forecasting," *W. Feltz and J. Murray, Eds., Remote Sensing Applications for Aviation Weather Hazard Detection and Decision Support. Proceedings of SPIE*, vol. 7088, p. paper 708805, August 2008.
- [13] N. Meinshausen, "Quantile regression forests," *Journal of Machine Learning Research*, vol. 7, pp. 983–999, 2006.
- [14] G. Biau, L. Devroye, and G. Lugosi, "Consistency of random forests and other averaging classifiers," *Journal of Machine Learning Research*, vol. 9, pp. 2015–2033, 2008.
- [15] J. K. Williams, R. Sharman, J. Craig, and G. Blackburn, "Remote detection and diagnosis of thunderstorm turbulence," in *Proceedings of SPIE*, vol. 7088. Remote Sensing Applications for Aviation Weather Hazard Detection and Decision Support, 2008.
- [16] U. Sharan and J. Neville, "Exploiting time-varying relationships in statistical relational models," in *Proceedings of the 1st SNA-KDD Workshop, 13th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2007.
- [17] —, "Temporal-relational classifiers for prediction in evolving domains," in *Proceedings of the IEEE International Conference on Data Mining*, 2008.
- [18] D. J. Allcroft, C. Glasbey, and M. Durban, "Modelling weather data," in *SCRI Annual Report 2001*, 2001, pp. 192–195.
- [19] C. Glasbey and D. J. Allcroft, "Spatio-temporal weather models," Talk at 39th French Statistical Association (SFDS) congress, Angers, June 2007, <http://www.bioss.sari.ac.uk/staff/chris/angers07.pdf>.
- [20] H. J. Miller and J. Han, Eds., *Geographic Data Mining and Knowledge Discovery*, 2nd ed. Chapman and Hall/CRC Press, 2009.
- [21] P. A. Valdes-Sosa, "Spatio-temporal autoregressive models defined over brain manifolds," *Neuroinformatics*, vol. 2, no. 2, pp. 239–250, 2004.
- [22] R. O. Weber and P. Talkner, "Some remarks on spatial correlation function models," *Monthly Weather Review*, vol. 121, no. 9, pp. 2611–2617, 1993.
- [23] R. Sharman, C. Tebaldi, G. Wiener, and J. Wolff, "An integrated approach to mid- and upper-level turbulence forecasting," *Weather and Forecasting*, vol. 21, pp. 268–287, 2006.
- [24] L. B. Cornman and B. Carmichael, "Varied research efforts are under way to find means of avoiding air turbulence," *ICAO Journal*, vol. 48, pp. 10–15, 1993.
- [25] T. Lane, R. Sharman, T. L. Clark, and H. Hsu, "An investigation of turbulence generation mechanisms above deep convection," *J. Atmos. Sci.*, vol. 60, no. 10, pp. 1297–1321, May 2003.
- [26] J. Williams, L. Cornman, J. Yee, S. Carson, G. Blackburn, and J. Craig, "NEXRAD detection of hazardous turbulence," in *44th American Institute of Aeronautics and Astronautics (AIAA) Conference on Aerospace Sciences*, 2006.
- [27] R. Frehlich and R. Sharman, "Estimates of turbulence from numerical weather prediction model output with applications to turbulence diagnosis and data assimilation," *Monthly Weather Review*, vol. 132, pp. 2308–2324, 2004.
- [28] D. D. Jensen and P. R. Cohen, "Multiple comparisons in induction algorithms," *Machine Learning*, vol. 38, no. 3, pp. 309–338, 2000.
- [29] K. Dwyer and R. Holte, "Decision tree instability and active learning," in *ECML '07: Proceedings of the 18th European conference on Machine Learning*. Berlin, Heidelberg: Springer-Verlag, 2007, pp. 128–139.