

1-1-2012

A Nonlinear Statistical Algorithm to Predict Daily Lightning in Mississippi

Erin Amanda Thead

Follow this and additional works at: <https://scholarsjunction.msstate.edu/td>

Recommended Citation

Thead, Erin Amanda, "A Nonlinear Statistical Algorithm to Predict Daily Lightning in Mississippi" (2012).
Theses and Dissertations. 214.

<https://scholarsjunction.msstate.edu/td/214>

This Graduate Thesis - Open Access is brought to you for free and open access by the Theses and Dissertations at Scholars Junction. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of Scholars Junction. For more information, please contact scholcomm@msstate.libanswers.com.

A nonlinear statistical algorithm to predict daily lightning in Mississippi

By

Erin Amanda Thead

A Thesis
Submitted to the Faculty of
Mississippi State University
in Partial Fulfillment of the Requirements
for the Degree of Master of Science
in Geosciences (Professional Meteorology)
in the Department of Geosciences

Mississippi State, Mississippi

December 2012

Copyright by
Erin Amanda Thead
2012

A nonlinear statistical algorithm to predict daily lightning in Mississippi

By

Erin Amanda Thead

Approved:

Andrew E. Mercer
Assistant Professor of Geosciences
(Director of Thesis)

Jamie L. Dyer
Associate Professor of Geosciences
(Committee Member)

Michael E. Brown
Associate Professor and Graduate
Coordinator of Geosciences
(Committee Member)

R. Gregory Dunaway
Professor and Interim Dean
College of Arts & Sciences

Name: Erin Amanda Thead

Date of Degree: December 15, 2012

Institution: Mississippi State University

Major Field: Geosciences (Professional Meteorology)

Major Professor: Dr. Andrew Mercer

Title of Study: A nonlinear statistical algorithm to predict daily lightning in Mississippi

Pages in Study: 57

Candidate for Degree of Master of Science

Recent improvements in numerical weather model resolution open the possibility of producing forecasts for lightning using indirect lightning threat indicators well in advance of an event. This research examines the feasibility of a statistical machine-learning algorithm known as a support vector machine (SVM) to provide a probabilistic lightning forecast for Mississippi at 9 km resolution up to one day in advance of a thunderstorm event. Although the results indicate that SVM forecasts are not consistently accurate with single-day lightning forecasts, the SVM performs skillfully on a data set consisting of many forecast days. It is plausible that errors by the numerical forecast model are responsible for the poorer performance of the SVM with individual forecasts. More research needs to be conducted into the possibility of using SVM for lightning prediction with input data sets from a variety of numerical weather models.

Key words: support vector machines, lightning, weather forecasting, statistical modeling, lightning forecasting, Mississippi

TABLE OF CONTENTS

LIST OF TABLES	iv
LIST OF FIGURES	v
CHAPTER	
I. INTRODUCTION	1
1.1 Overview of the Project	1
1.2 Literature Review.....	2
1.2.1 Review of Lightning Forecasting Parameterizations and Techniques	2
1.2.2 Review of Support Vector Machines.....	7
1.3 Research Objectives.....	8
II. DATA AND METHODOLOGY.....	10
2.1 Data Description	10
2.2 Methodology.....	11
2.2.1 WRF setup	13
2.2.2 SVM overview	17
2.2.3 Lightning data	20
2.2.4 Interpolation of data.....	20
2.2.5 Selection of parameters.....	22
2.2.6 Sampling of data	25
2.2.7 Validation of the SVM.....	26
2.2.8 Methods summary.....	27
III. RESULTS	28
3.1 Prediction of a large data set.....	28
3.2 Prediction of individual cases	29
3.2.1 3 October 2002.....	29
3.2.2 24 November 2001.....	32
3.2.3 20 October 2004.....	35
3.2.4 22 September 2006	38
3.2.5 16 September 2005	41
3.3 Interpretation of results	44
3.3.1 NARR and WRF-ARW limitations	44

3.3.2	Rarity of the event.....	46
IV.	SUMMARY AND CONCLUSIONS	47
	REFERENCES	50
	APPENDIX	
A.	CASES USED IN THE RESEARCH.....	55

LIST OF TABLES

2.1	WRF physics parameterizations	16
2.2	Parameters subjected to permutation testing.....	24
3.1	Contingency statistics for full predictor set	28
3.2	Contingency statistics for 3 October 2002.....	32
3.3	Contingency statistics for 24 November 2001	35
3.4	Contingency statistics for 20 October 2004.....	38
3.5	Contingency statistics for 22 September 2006.....	41
3.6	Contingency statistics for 16 September 2005.....	44
A.1	Lightning cases and total lightning strikes in Mississippi	56

LIST OF FIGURES

2.1	Histogram of strike counts in all autumn 2001-2007 lightning days.....	12
2.2	Histogram of strike counts in 90 selected lightning days	13
2.3	Domain of the WRF simulation.....	15
2.4	Support vector machine function.....	18
2.5	SVM kernelization.....	19
2.6	Distribution of lightning day counts in 9 km grid.....	21
3.1	Observed lightning on 3 October 2002.....	30
3.2	SVM-predicted lightning on 3 October 2002	31
3.3	Observed lightning on 24 November 2001	33
3.4	SVM-predicted lightning on 24 November 2001	34
3.5	Observed lightning on 20 October 2004.....	36
3.6	SVM-predicted lightning on 20 October 2004	37
3.7	Observed lightning on 22 September 2006.....	39
3.8	SVM-predicted lightning on 22 September 2006	40
3.9	Observed lightning on 16 September 2005.....	42
3.10	SVM-predicted lightning on 16 September 2005	43

CHAPTER I

INTRODUCTION

1.1 Overview of the Project

Lightning is one of the leading causes of severe weather-related fatalities in the United States. However, localized operational forecasts for cloud-to-ground (CG) lightning are not currently issued. Lightning of any kind is difficult to predict in advance, so regional thunderstorm activity is often used as a proxy for CG strikes. Lightning activity in storms that produce high rates of CG lightning is only forecasted in the very short range (three hours or less) or tracked if it is already occurring. That is, if any CG lightning activity is mentioned, it is mentioned only in severe thunderstorm warnings that are issued based on dangerous phenomena in the storm or in “special weather statements” that are issued at the discretion of the local National Weather Service office. The only regular indicator in these products that lightning is occurring is the fact that they are issued for thunderstorms, thereby implying by definition the possibility of lightning.

The lack of local operational lightning forecasts and lightning products is caused by the difficulty in explicitly forecasting CG lightning threats. Lightning prediction requires high-resolution data for reasonably accurate forecasts to be expected, and until recently, numerical weather model prediction has not been accurate enough in forecasting storm-scale features to justify the production of high-resolution forecasts. Improvement in model performance owing to greater computational power opens the possibility of

producing forecast products for CG lightning using indirect lightning threat indicators well in advance of an event. Such forecasts would have to be statistically based, as numerical models are incapable of explicitly resolving lightning within a thunderstorm using storm dynamics. This research examines a statistical learning-algorithmic approach for forecasting lightning probabilistically in high resolution, up to 1 day in advance of an anticipated event.

1.2 Literature Review

1.2.1 Review of Lightning Forecasting Parameterizations and Techniques

The electrification of cumulonimbi is believed to be associated with charge separation occurring in the process of ice formation. Two primary processes have been documented to produce electrification in clouds (Saunders 1992)—non-inductive and inductive charging.

Inductive charging involves the collision of ice fragments with supercooled water, which results in the formation of graupel. In the presence of an existing electric field, electrostatic induction creates a positive charge in the ice crystals and a negative charge in the graupel, and updrafts separate the two forms of frozen water, lofting the ice vertically, while graupel falls downward. This process, however well-understood, is not believed to be able to account for the levels of charging observed in thunderstorms (Jennings 1975). It also has difficulty accounting for the levels of charging observed in particles in early storm development (Saunders 1992), when a strong electric field has not yet been established.

The non-inductive mechanism also requires ice, graupel, and supercooled water, but it does not require an existing electric field that produces induction. Laboratory

studies simulating conditions within developing thunderstorms (Reynolds et al. 1957, Takahashi 1978) have found that charge separation occurs in association with the breaking of graupel particles in updrafts and the contact of these rime pieces with each other. The presence of supercooled water in the surrounding environment increases the amount of charge acquired by the rime pellets (Jayaratne and Griggs 1991). The sign of the charge is influenced by the phase change that graupel particles undergo as they break apart (Williams et al. 1991, Hallett and Saunders 1979); particles undergoing deposition generally charge positively and particles undergoing sublimation generally charge negatively.

In contrast with inductive charging, researchers have found that non-inductive charging accounts for the observed amounts of electrification in thunderstorms very well. Takahashi (1978) observed this result in laboratory simulations of intracloud conditions, and Fierro et al. (2008) found that non-inductive charging accounted for the majority of charges in a tropical squall line.

High numbers of ice particle collisions result in higher levels of charge, indicating that the amount of ice in a thunderstorm can be used as a proxy for the electrical charging potential. Models are not able to resolve microscale processes such as the disintegration of graupel, which occurs on the scale of millimeters or less. Due to this resolution issue, schemes have been formulated that use regression methods to predict lightning formation, based on variables with high correlation to lightning activity.

One scheme was devised by Kitzmiller et al. (2000) of the National Weather Service. This scheme is based upon a statistical regression of several predictor variables for both CG lightning prediction and rainfall rate. These predictors were chosen from a

larger set of candidates. Probability values were calculated for each candidate parameter based on its ability to predict lightning strikes in a 40 km grid in 15 minute periods. This scheme used different moisture and instability parameters, including mean relative humidity, precipitable water, K index, 850 mb lifted index, moisture divergence, and 6-hour precipitation from Nested Grid Model (NGM) forecast grids. Additionally, the regression model utilized observed radar parameters representing the summation of modeled grid boxes containing high-level dBZ echoes and satellite-measured infrared temperature as indicators of the precipitation intensity and cloud top temperature of thunderstorms, respectively. The model also used a predictor representing the CG lightning strike rate in a given grid box over the past 15 minutes. Their lightning observation data were provided by the Marshall Space Flight Center. The radar, infrared satellite, and lightning strike parameters were extrapolated forward 3 hours using NGM 700-500 mb wind vector data, and the regression analysis was performed on the extrapolated results.

They found that, of the variables tested, the best predictor for CG lightning was the radar parameter, which outperformed even the observed lightning strike rate. They speculated that the radar parameter was more conservative and less time-sensitive than the lightning strike rate, implying overfitting with the lightning strike parameter. However, the extrapolation of radar data is error-prone and imprecise in general. Storm-scale atmospheric conditions ahead of radar-indicated storm systems may change rapidly. Furthermore, the use of observed lightning data in a lightning forecast necessitates a short lead time, and indeed, this forecasting algorithm was intended only for “short-range”

forecasts at most 3 hours in advance of an event. For their final lightning forecast algorithm, the probability of detection was 0.67 and the false alarm rate was 0.53.

Mazany et al. (2002) developed a different scheme for forecasting lightning. Rather than producing forecasts for an entire area over a period of time, their intent was to provide a means of forecasting a developing thunderstorm's first strike 90 minutes in advance, and to that end they used a parameter based on integrated precipitable water vapor as measured by global positioning system (GPS) satellites. They used a logistic regression to determine which of 23 variables had the best correlation for lightning forecasts and found that GPS integrated precipitable water vapor performed the best. The test used backward selection to progressively eliminate possible predictors according to their statistical significance. Mazany et al. (2002) focused on lightning events in the state of Florida, as their objective was to provide a forecasting tool for use by the National Aeronautics and Space Administration (NASA) in determining whether to launch the space shuttle.

Bright et al. (2005) developed a parameter that they call the "cloud physics thunder parameter" for making forecasts about the formation of thunderstorms and therefore lightning. Their formula used Convective Available Potential Energy (CAPE) from the 0°C to -20°C levels of the atmosphere, which is generally the most unstable region, based on the idea that some amount of uplift is required to lift graupel above an area of the storm detrimental for electrification. Jayaratne et al. (1983) noted that between -15°C and -20°C, the sign of charged graupel reverses, resulting in a significant reduction of lightning activity in this narrow vertical region. They also used the equilibrium level temperature in their formula and determined that this value needs to be

less than or equal to -20°C so that the top of the cloud can also extend past the charge-reversal zone.

Bright et al. (2005) further modified this parameter with two constants, one of which (a constant, K , defined as 100 J kg^{-1}) was determined strictly by experimentation. In creating a hypothetical operational forecast, they modified their technique further by multiplying their computed cloud physics thunder parameter by the probability of precipitation greater than or equal to 0.01” (as computed by the Short-Range Ensemble Forecast suite of weather models). They found good results with their cloud physics thunder parameter, with their forecasts showing 10 to 15 percent improvement over climatology as calculated in the Brier Skill Score (Wilks 2006). The cloud physics thunder parameter is currently calculated in post-processing by the Short-Range Ensemble Forecast (SREF) system.

Yet another approach was employed by Fierro et al. (2006) in analyzing a particular event in which an outbreak of discrete supercells occurred. The storms in their research occurred on 2 June 1995 over Texas, and they initialized the models with an idealized horizontal environment typical of an outflow boundary. They focused on the effects upon a supercell of crossing this boundary to the cool side. They found that simulated supercells that crossed the boundary intensified rapidly, whereas those that did not cross the boundary did not. Intensification of a supercell was associated with increased electrification as the updraft strength increased, as measured by observed lightning strikes from the National Lightning Detection Network (NLDN) database (Orville 2008). In addition, they found that graupel amounts, hail amounts, and maximum radar reflectivity increased when a simulated storm crossed the outflow

boundary. Their research indicates a relationship between thermal characteristics of the surrounding mesoscale environment and cloud parameters associated with ice formation.

1.2.2 Review of Support Vector Machines

The statistical learning technique that was used in this research, support vector machines (SVM) (Burges 1998, Hearst et al. 1998, Cristianini and Shawe-Taylor 2000), is a type of learning machine well suited for classification. SVMs use a technique known as decision hyperplanes. Similar to a decision line, a decision hyperplane can be extended to n -dimensional space (hyperspace). These planes demarcate the outcome data in a binary fashion, with all points on one side corresponding to one outcome and all on the other side to the other outcome.

Support vector machines have a computational complexity of $O(n^2)$ to $O(n^3)$, depending upon the amount of optimization required. This is a polynomial-time algorithm, which is to say that its computation time grows “manageably” (as opposed to exponential-time algorithms) as the problem size increases. Due to their computational tractability, support vector machines and other learning algorithms have been used in meteorological research.

One such machine learning algorithm, the neural network, has been used in classification and forecasting a variety of events, such as automated cloud observation (Aviolat et al. 1998), estimation of rainfall from radar (Liu et al. 2001), tornado prediction from radar signature (Marzban and Stumpf, 1996), and coastal water level prediction from weather station data (Han and Shi 2008).

Despite the advances made with neural networks, it has been shown that support vector machines are less prone to overfitting (Burges 1998), an important consideration

for meteorological data sets and forecasting. Mercer et al. (2009) have used support vector machines to classify severe weather events as tornadic or nontornadic given predictors corresponding to various physical parameters, with the SVM approach showing skill over a logistic regression approach. The technique has also been used in global cloud mask algorithms (Garay et al. 2003), aerosol modeling (Ackerman et al., 2004), classification of satellite radiance data into cloud types (Lee et al. 2003), and downslope windstorm forecasting in Colorado (Mercer et al. 2008).

1.3 Research Objectives

The research question to be answered was whether a statistical algorithm could be formulated that could skillfully predict the probability of CG lightning in a 9 km area up to 1 day in advance. Since the ultimate goal was to produce a model that could be used for operational lightning forecasting, any such algorithm would have to be computationally tractable and produce its forecast quickly enough for operational forecasting usage.

As the work of Kitzmiller et al. (2000), Mazany et al. (2002), and Bright et al. (2005) has indicated, it is possible to formulate a skilled model based strictly on statistical correlation without respect to any cloud-physical process. This research determined an optimal combination of parameters to predict the daily CG lightning threat. After this determination was made, the research utilized support vector machines (SVMs) to formulate a CG probability algorithm for Mississippi to determine SVM's effectiveness at predicting daily CG lightning activity.

The remaining sections of this document contain specific details about the data sets and methods that were used in this research, as well as a description of the results

obtained and a summary of the conclusions. Section 2, Data and Methodology, contains a detailed description of data sources and modeling tools, including background information on the data and tools, resolution, domains, known sources of error, limitations, and configuration details specific to this research. Section 3, Results, contains qualitative (graphical) and quantitative (statistical) analysis of the results of the research. This section also contains an explanation of factors relating to the data sets, the modeling tools, and the nature of the problem that may have affected the results. Section 4, Summary and Conclusions, is a summary of the research with an emphasis on the results obtained. This section also contains a subsection proposing possible related topics for future research.

CHAPTER II

DATA AND METHODOLOGY

2.1 Data Description

The predictor data given as input to the SVM were derived from the Weather Research and Forecasting-Advanced Research WRF (WRF-ARW) model, version 3.2.1. (Skamarock et al. 2005). The WRF model was initialized with the North America Regional Reanalysis (NARR) data (Mesinger et al. 2006). These data have 32 km grid spacing with 29 vertical levels and 3-hour temporal resolution. The NARR data were reanalyzed from recorded weather observations using the North American Mesoscale (NAM) model, formerly known as the Eta model, and assimilated with the Eta data assimilation system. Two known limitations of the data were that the Gulf of California low level jet is too strong in the summer, and surface wind stress is insufficiently precise (Ebisuzaki and Rutledge 2004), neither of which was relevant for this research.

An additional limitation of the NARR data, as described by Mesinger et al. (2006), was that 2-meter temperature fields in these data do not correspond well to 2-meter temperatures fitted to the observations generated by tropospheric rawinsondes. Diurnal variation of land temperatures owing to boundary layer processes results in a large impact on lower-tropospheric temperatures and winds, owing to the NAM's inability to limit vertical influence of these parameters. To compensate for this problem, 2-meter temperatures in the NARR data show less diurnal variation than actually occurs.

This limitation of the NARR data was relevant to this research, as temperature at all levels was used in the WRF simulation.

In addition to the WRF output, the SVM trained with the observed cloud-to-ground (CG) lightning strike data as taken from the National Lightning Detection Network (NLDN) database. These data are recorded by a network of sensors that measure surges in electromagnetic radiation. The readings are sent to a system owned by Vaisala, Inc., from which point they may be disseminated to organizations that subscribe to the lightning data. The data are also archived in the NLDN database for research usage. The observations record the date and time (to the nearest millisecond), latitude/longitude location (to thousandths of a degree), signal strength (in kA), polarity, and number of return strokes (Orville 2008). The NLDN sensors are able to detect flashes above 5 kA at 80 to 90 percent (Cummins et al. 1998, Burrows et al. 2005), with a location error ranging from between 435 m and 625 m (Idone et al. 1998). This study, however, included all strikes measured by the sensors, including those registering at lower levels of electrical current.

2.2 Methodology

The formulation of the SVM algorithm required three primary steps: generation of the predictor data for the SVM using WRF-ARW simulations, training of the SVM, and validation of the SVM's forecasts.

The research used 90 cases from 2001-2007 in which CG lightning occurred in Mississippi during a 24-hour period (0000 UTC to 2359 UTC). All cases were randomly sampled from meteorological autumn, defined here as the months of September, October, and November. The 90 cases were subdivided into quintiles based upon the number of

strikes reported and randomly selected the same number of cases for each quintile. The purpose of this division was to ensure that the 90 randomly selected cases were representative of the full (autumn 2001-2007) set of days in terms of lightning strike counts; i.e., that high-impact days, for example, were not over-sampled among the 90 cases. Histograms of the full set (Figure 2.1) and of the set of 90 (Figure 2.2) are comparable, indicating that the sampling did reflect the distribution of the full data set.

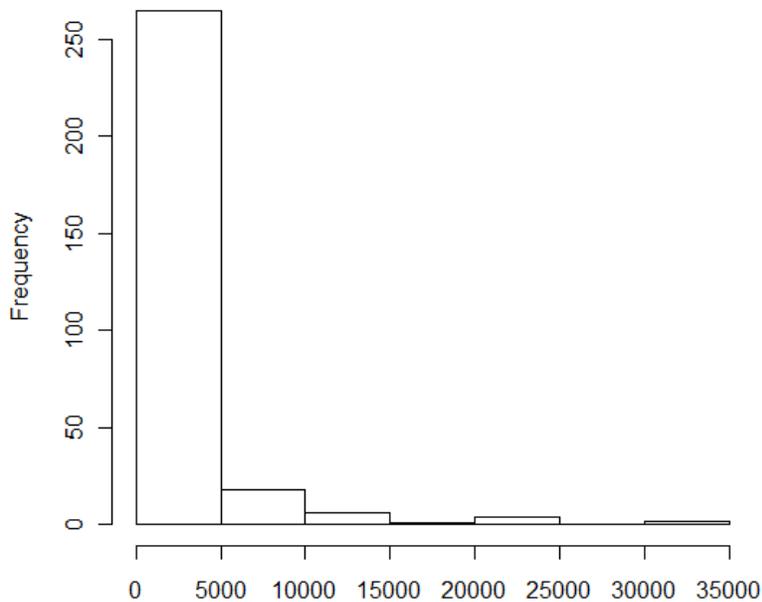


Figure 2.1 Histogram of strike counts in all autumn 2001-2007 lightning days

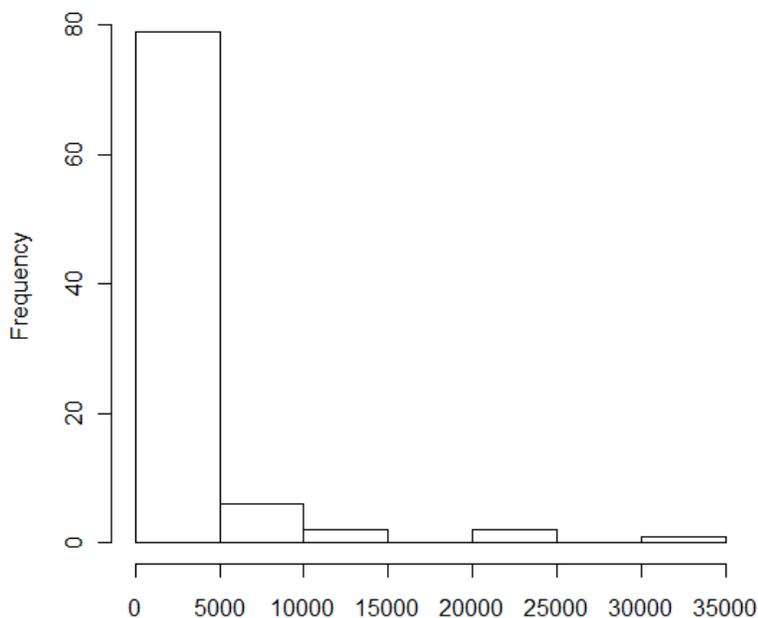


Figure 2.2 Histogram of strike counts in 90 selected lightning days

Each case was simulated in the WRF-ARW model 30 hours prior to 0000 UTC on the event day to 6 hours after 0000 UTC. For example, the case for September 2, 2001 included 1800 UTC from September 1 through 0600 UTC from September 3. This simulation range was chosen to ensure the ability to forecast 12 hours prior to the event day (which began at 0600 UTC, or local midnight) and account for all lightning activity on that day as observed in local time.

2.2.1 WRF setup

The NARR data are not prognostic tools, but instead represent reanalyses of past meteorological days. Thus, they cannot be used in making operational forecasts. Support vector machines should be trained with the same type of data that they will use operationally, so it was deemed necessary to use data from a numerical weather model for training the SVM. The WRF-ARW weather model, which was developed to be

readily configurable for specific research needs (Skamarock and Klemp 2007), was used for this research. The reason for this choice is that the WRF core is used in the operational North American Mesoscale (NAM) model, the primary finite-element non-hydrostatic mesoscale model run by the National Centers for Environmental Prediction (NCEP). This model (and other models that use NAM model output as their input) contains the same biases and limitations of the WRF. An SVM-based operational lightning forecast would use NAM-derived model data as predictors for the SVM; therefore it was necessary to train the SVM with this type of data.

The WRF-ARW model configuration that was chosen used a single domain representing a grid box enclosing the state of Mississippi and additional area for accommodating rapidly changing boundary conditions (Figure 2.3). The time step of the configuration was 15 seconds, with output files generated for each hour, and the spatial resolution was 3 km. WRF-ARW outputs to 40 vertical levels, and this vertical output was then interpolated in post-processing to every 25 mb atmospheric pressure level from 1000 mb to 100 mb (the top of the NARR vertical domain).

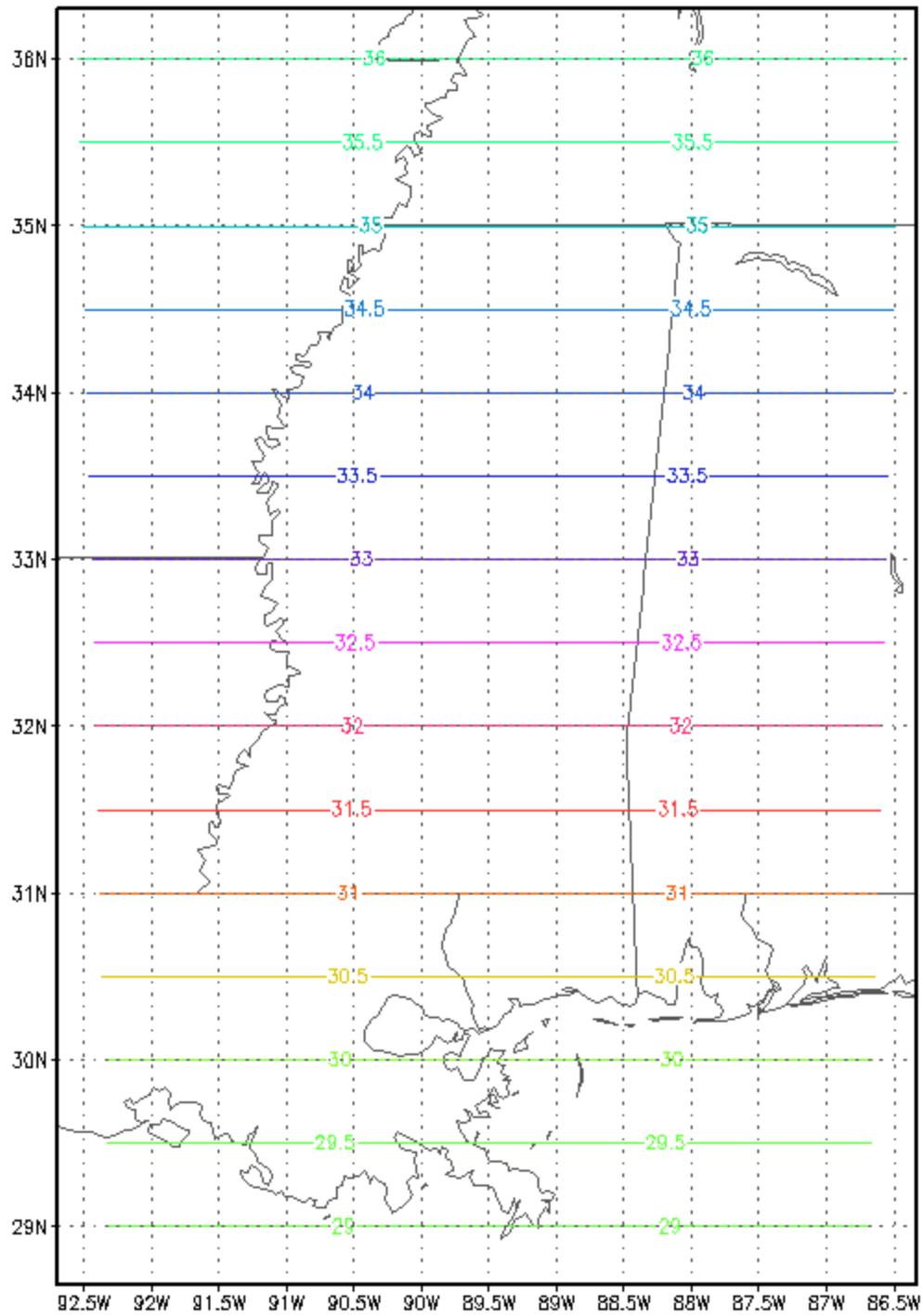


Figure 2.3 Domain of the WRF simulation

The Thompson et al. cloud microphysics scheme (Hall et al. 2005) was used. This scheme was developed for ice formations in high-resolution simulations (Thompson et al. 2006) and therefore is well-suited to this research. Since the WRF model was run at a high resolution, a cumulus physics scheme for parameterizing simulated clouds was not used.

The suite of physics schemes used for this research are listed in Table 2.1.

Table 2.1 WRF physics parameterizations

WRF physics option	Configuration	Reference
Cloud microphysics	Thompson et al.	Thompson et al. 2006
Longwave radiation	Rapid Radiative Transfer Model	Mlawer et al. 1997
Shortwave radiation	Dudhia	Dudhia 1989
Surface layer	MM5-derived	Dudhia 1996
Land surface	5-layer thermal diffusion	Dudhia 1996
Urban surface	None	
Planetary boundary layer	Yonsei University	Hong and Pan 1996
Cumulus physics	None	

Koch et al. (2005) described in detail the applicability of WRF simulated radar reflectivity to analyzing mesoscale and storm-scale weather phenomena. They observed that in convective thunderstorm events, high-resolution (2 km horizontal grid) WRF could model storm-scale structure in its simulated reflectivity product. Moreover, they found that the WRF-ARW modeled strong echoes (>50 dBZ) better than the WRF-Nonhydrostatic Mesoscale Model (WRF-NMM), which is used by NCEP for operational forecasts. The reason for this is directly related to the modeling of water and ice within clouds; the WRF-NMM simplifies its computations by assuming a maximum concentration of precipitated ice, a mean size of 1 mm for a precipitating ice crystal, and a fixed raindrop size (0.45 mm) for rain over 1 g/m^3 . Although the WRF-NMM core is

used in operational models rather than the WRF-ARW, the differences between the two are differences of precision in microscale physics. The general biases of the model remain the same. Furthermore, the product that this research produced is a probabilistic forecast for 9 km grid spaces. The use of probability in the final product minimized any possible overfitting that might have occurred as a result of using training input that is more precisely computed in microscale than the input that forecasters would use operationally.

The raw WRF output was post-processed using the ARWpost visualization software, and the parameters given to the SVM and logistical regression model were taken from the ARWpost output files. The Grid Analysis and Display System (GrADS) and a Fortran binary-to-text conversion program were chosen to convert ARWpost files to a human-readable text format suitable for statistical analysis. It is this format that was analyzed in the SVM and logistical regression models and used to generate a lightning forecast.

2.2.2 SVM overview

SVMs use quadratic programming optimization to determine the best location of the decision hyperplane, as an infinite number of possible hyperplanes will exist for any given data set.

The equation of a hyperplane can be generalized as

$$w^T x + b = 0 \tag{2.1}$$

where w is a vector of weights, x is a vector of covariates, and b is an intercept. The classes of the data points are found by reformulating this equation as an inequality.

Points on the side where the value is positive represent one class, and points on the negative side represent the other (Cristianini and Shawe-Taylor 2000) (Figures 2.4 and 2.5).

The quadratic optimization problem for a hyperplane is given as

$$\max F(\Lambda) = \sum_{i=1}^l \lambda_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \lambda_i \lambda_j y_i y_j x_i x_j \quad (2.2)$$

subject to $\sum_{i=1}^l \lambda_i y_i = 0$ $\lambda_i \geq 0$, where the values of λ_i and λ_j represent Lagrange multipliers, x_i and x_j represent the covariates, y_i and y_j are solutions to the hyperplane equation.

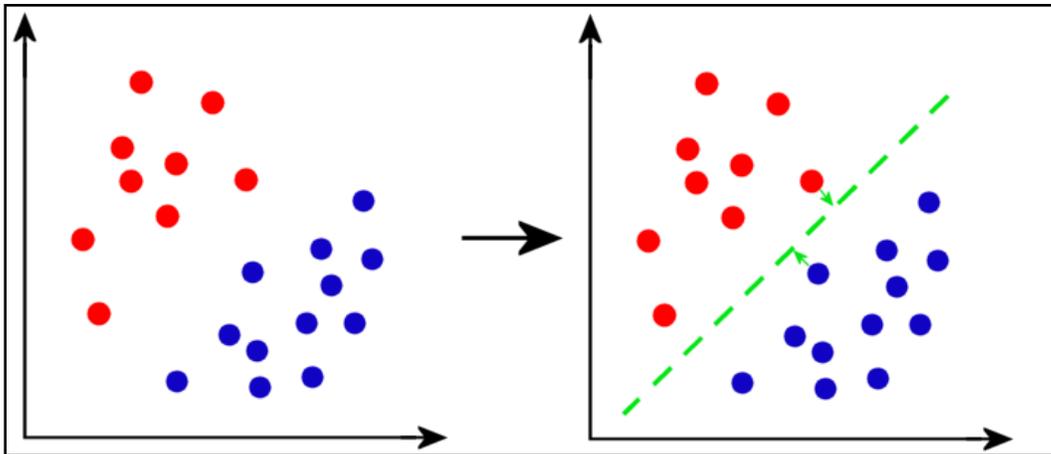


Figure 2.4 Support vector machine function

Blue points represent one class of data and red points represent another. The dashed line represents a decision hyperplane. The dashed diagonal line in the right image is the decision hyperplane. The arrows in the right image are the “support vectors” of the algorithm; the points that these arrows are on are closest to the hyperplane and it is the distance between them that must be maximized.

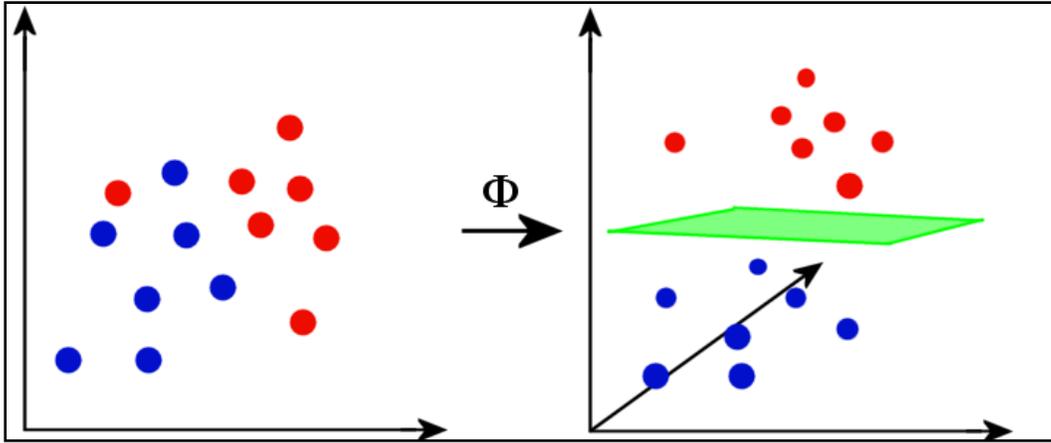


Figure 2.5 SVM kernelization

This image depicts a data sample for which no 2-D hyperplane exists and its kernelization into a higher dimension via a kernel function Φ .

The SVM training used the 10 predictors selected by permutation testing. The values of these predictors were taken from the lightning data set resampled as described in §2.2.5. Numerous configurations of SVM parameters were tested. The kernel functions tested were as follows:

1. A radial basis kernel function

$$k(x, y) = e^{(-1/2\sigma^2)\|x-y\|^2} \quad (2.3)$$

2. A polynomial kernel function

$$k(x, y) = (x^T y + 1)^p \quad (2.4)$$

Different costs were also tested. The cost is a parameter that controls the sensitivity of the SVM to the training data. A higher cost value increases the cost of training errors and results in a more accurate model, but risks overfitting.

Little difference was found in the predictive power of the SVM across the varied configurations. The SVM configuration that provided the results described in §3 of this document was trained with a cost of 1000 and a radial kernel.

2.2.3 Lightning data

The SVM is programmed to distinguish between distinct classifications of data. For this research, it required binary (yes/no) data for its predictand. Therefore, lightning strikes in the state of Mississippi were extracted from the NLDN strike database for the 90 cases and a text file of ones (representing yeses) and zeroes (representing nos) for each 3 km grid space in Mississippi for each of the events was created. The lightning strike data set matched the predictor data set in grid size and geographical location.

2.2.4 Interpolation of data

Once the parameters were selected, the parameter and lightning data were interpolated from a 3 km by 3 km grid scale to 9 km by 9 km. This was because an initial run of the SVM with the 3 km data proved to be intractable. The data set included 102,750 data points per predictor per case, and with 90 cases and 10 predictors for each case, this resulted in over 92 million data points. The length of time that the SVM took to run in this configuration was deemed infeasible for a potential operational forecasting product. After 24 hours of continuous calculation, the SVM was still running and had not converged on a solution. Clearly this length of time is not viable for an operational forecasting product. The predictor and lightning data were therefore interpolated to 9 km. The interpolation of the lightning data retained the binary character of the original 3 km lightning strike data file; it did not accumulate lightning strikes in the larger grid.

The interpolation resulted in a grid of 11,500 data points per predictor per case in Mississippi.

It was observed that, even with a slightly coarser resolution of 9 km, lightning days were still a rare occurrence. The highest number of days that any particular grid box experienced lightning was 13, out of a possible 90. Almost 50 percent of the squares never reported lightning even though it occurred elsewhere in the state. The distribution of the total number of lightning days in each 9 km square over all 90 cases was calculated. This distribution is shown in Figure 2.6.

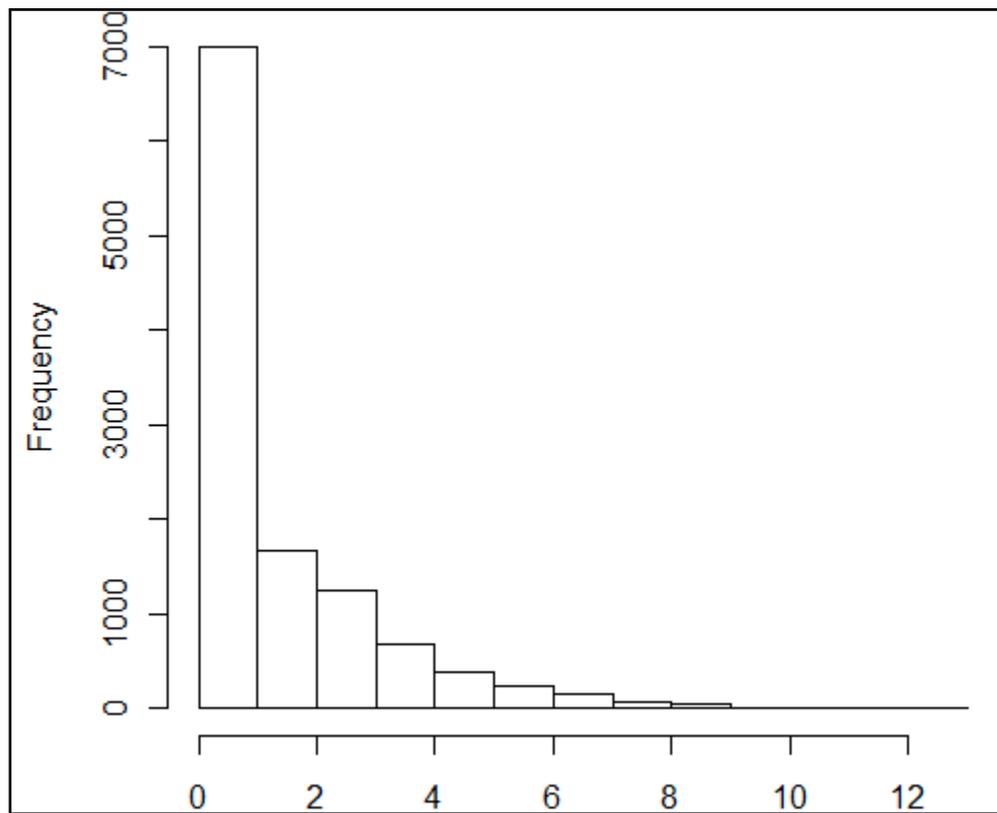


Figure 2.6 Distribution of lightning day counts in 9 km grid

2.2.5 Selection of parameters

Permutation testing was used to determine the final predictors (Efron and Tibshirani 1993) given to the SVM. This statistical technique resamples data to determine if the means of two distributions are statistically significantly different. The test resamples a given number of times (the number of permutations), calculates the means of the two distributions, compares these values to the mean of the full, pooled data set. The number of times that the pooled mean difference is greater than the individual mean difference is divided by the number of permutations. This is the p-value. The permutation testing used in this research used 2000 permutations to calculate the p-values.

Permutation testing was chosen over other tests, such as the t-test, because it does not require a normalized data set (unlike the t-test), and a rare-events data set such as lightning cannot be assumed to be normally distributed.

Since the scope of this project is individual, pointwise lightning prediction, it was deemed necessary to choose parameters by a method that is sensitive to the actual occurrence of lightning in a grid box. The permutation test used parameter values only from the 9 km grid boxes in which, over the 90 lightning days, 5 or more lightning days occurred. As Figure 2.6 shows, lightning was very rare in any particular 9 km space. This sampling was to give the permutation test a data set containing a larger percentage of values corresponding to lightning events. The groups that were compared were the parameter values for the lightning days and the parameter values for non-lightning days within the larger set of 90 cases where lightning occurred elsewhere in the state.

Four values of the parameters to be tested were extracted from each of the 90 cases: the values for 0600Z, 1200Z, 1800Z, and 00Z on the following day. One value of each parameter was extracted for every grid box in the domain. A p-value was calculated for each parameter. If a p-value was less than the rejection threshold, which was calculated using bootstrap confidence intervals of 95%, the hypothesis that the means of this value for lightning events and non-lightning events are the same was rejected, and the parameter was used in training the SVM. A variety of parameters were tested, as listed in Table 2.2. These parameters were associated with convective activity and cloud electrification as lightning predictors. The final predictors chosen were the best performing parameters.

The permutation tests did not reveal any statistically significant differences between parameter values for lightning events and parameter values for non-lightning events in a grid box. Therefore, the SVM predictor values were chosen according to the parameters that had the lowest confidence interval; i.e., the closest to statistical significance, and whose p-values were statistically significantly lower than those for other predictors. Predictors chosen were the planetary boundary layer height at 06Z and 12Z, accumulated upward surface heat flux at all time steps (06Z, 12Z, 18Z, and 00Z the following day), 500 mb north-south wind at 12Z and 18Z, and 250 mb north-south wind at 12Z and 18Z.

Table 2.2 shows a list of parameters that were tested, the results of each permutation test, and the bootstrap confidence intervals for the permutation test that used only grid boxes with 5 or more lightning days.

Table 2.2 Parameters subjected to permutation testing

Parameter	Atmospheric level	95% Confidence Interval				Reject/Keep
		06Z	12Z	18Z	00Z	
Perturbation dry air mass in column (mu)	1000 mb	0.33	0.31	0.37	0.41	Reject
Base state dry air mass in column (mub)	1000 mb	1.00	1.00	1.00	1.00	Reject
Water vapor mixing ratio at 2 m (q2)	1000 mb	0.47	0.46	0.40	0.37	Reject
Temperature at 2 m (t2)	1000 mb	0.51	0.46	0.43	0.37	Reject
Skin sea surface temperature	1000 mb	0.51	0.46	0.40	0.37	Reject
Accumulated total cumulus precipitation (rainc)	1000 mb	1.00	1.00	1.00	1.00	Reject
Accumulated total grid scale precipitation (rainnc)	1000 mb	0.48	0.41	0.41	0.40	Reject
Accumulated total grid scale graupel (graupelnc)	1000 mb	0.48	0.46	0.43	0.42	Reject
Accumulated total grid scale hail (hailnc)	1000 mb	1.00	1.00	1.00	1.00	Reject
Planetary boundary layer height (pblh)	1000 mb	0.26	0.29	0.36	0.52	Keep 06Z Keep 12Z
Upward heat flux at the surface (hfx)	1000 mb	0.36	0.40	0.35	0.49	Reject
Upward moisture flux at the surface (qfx)	1000 mb	0.43	0.44	0.43	0.53	Reject
Latent heat flux at the surface (lh)	1000 mb	0.43	0.44	0.43	0.53	Reject
Accumulated upward heat flux at the surface (achfx)	1000 mb	0.30	0.30	0.25	0.25	Keep 06Z Keep 12Z Keep 18Z Keep 00Z next day
Accumulated upward latent heat flux at the surface (aclhf)	1000 mb	0.50	0.50	0.50	0.47	Reject
Dew point temperature at 2 m (td2)	1000 mb	0.48	0.46	0.42	0.37	Reject
Relative humidity at 2 mb (rh2)	1000 mb	0.42	0.60	0.31	0.47	Reject
Sea level pressure (slp)	1000 mb	0.37	0.35	0.38	0.43	Reject
Maximum simulated reflectivity (max_dbz)	1000 mb	0.45	0.35	0.49	0.47	Reject

Table 2.2 (continued)

Maximum Convective Available Potential Energy (mcape)	1000 mb	0.37	0.37	0.39	0.43	Reject
Lifted Condensation Level (lcl)	1000 mb	0.37	0.46	0.44	0.49	Reject
Level of Free Convection (lfc)	1000 mb	0.37	0.43	0.39	0.42	Reject
East-west wind at 10 mb (u10)	1000 mb	0.41	0.39	0.38	0.43	Reject
North-south wind at 10 m (v10)	1000 mb	0.31	0.35	0.42	0.48	Reject
Water vapor mixing ratio (qvapor)	850 mb	0.45	0.40	0.44	0.40	Reject
Cloud water mixing ratio (qcloud)	850 mb	0.42	0.43	0.43	0.48	Reject
Rain water mixing ratio (qrain)	800 mb	0.45	0.38	0.45	0.48	Reject
Dew point temperature (td)	700 mb	0.40	0.41	0.43	0.47	Reject
Relative humidity (rh)	700 mb	0.40	0.36	0.42	0.48	Reject
East-west wind (u)	500 mb	0.45	0.46	0.44	0.32	Reject
North-south wind (v)	500 mb	0.35	0.25	0.23	0.31	Keep 12Z Keep 18Z
Geopotential (geopt)	500 mb	0.44	0.43	0.40	0.38	Reject
Geopotential height (height)	500 mb	0.44	0.43	0.40	0.39	Reject
Vertical wind (w)	450 mb	0.43	0.43	0.45	0.49	Reject
Air temperature (tk)	450 mb	0.48	0.52	0.49	0.48	Reject
Reflectivity (dbz)	450 mb	0.41	0.37	0.45	0.50	Reject
Convective Available Potential Energy (cape)	450 mb	0.68	0.74	0.70	0.68	Reject
Graupel mixing ratio (qgraup)	300 mb	0.42	0.41	0.42	0.43	Reject
East-west wind (u)	250 mb	0.53	0.50	0.52	0.45	Reject
North-south wind (v)	250 mb	0.36	0.30	0.28	0.39	Keep 12Z Keep 18Z

2.2.6 Sampling of data

A predictor data set was created using data points from the 90 cases in which lightning occurred in the state. Only 9 km boxes in which lightning occurred were used.

All the squares that reported lightning were included for the yes points, and two times as many data points (from the same grid boxes) for events in which lightning did not occur in that grid box were used for the no points. This therefore resulted in a set of 33 percent yes events and 66 percent no events, over 52,158 total data points per predictor. The decision to use only days in which lightning occurred somewhere in Mississippi was justified for a prospective operational forecast product on the basis of forecaster awareness of potential thunderstorm days; it is assumed that forecasters would not need a statistical forecast model to predict lightning on days when conditions are prohibitive or highly unfavorable for thunderstorms (i.e. high pressure scenarios). The decision to oversample yes events (with respect to the actual, much lower percentage of yes events over the 9 km grid) was justified on the basis of providing the SVM with enough yes events to distinguish them in its forecasts. Initial runs of the SVM with yes-event percentages closer, or identical, or the actual percent coverage of lightning resulted in SVMs that were unable to distinguish between the yes and no events at all.

Because the predictors included a variety of meteorological parameters and the values thereof exhibited a wide range, it was deemed necessary to normalize the predictor data set before supplying it as input to the SVM.

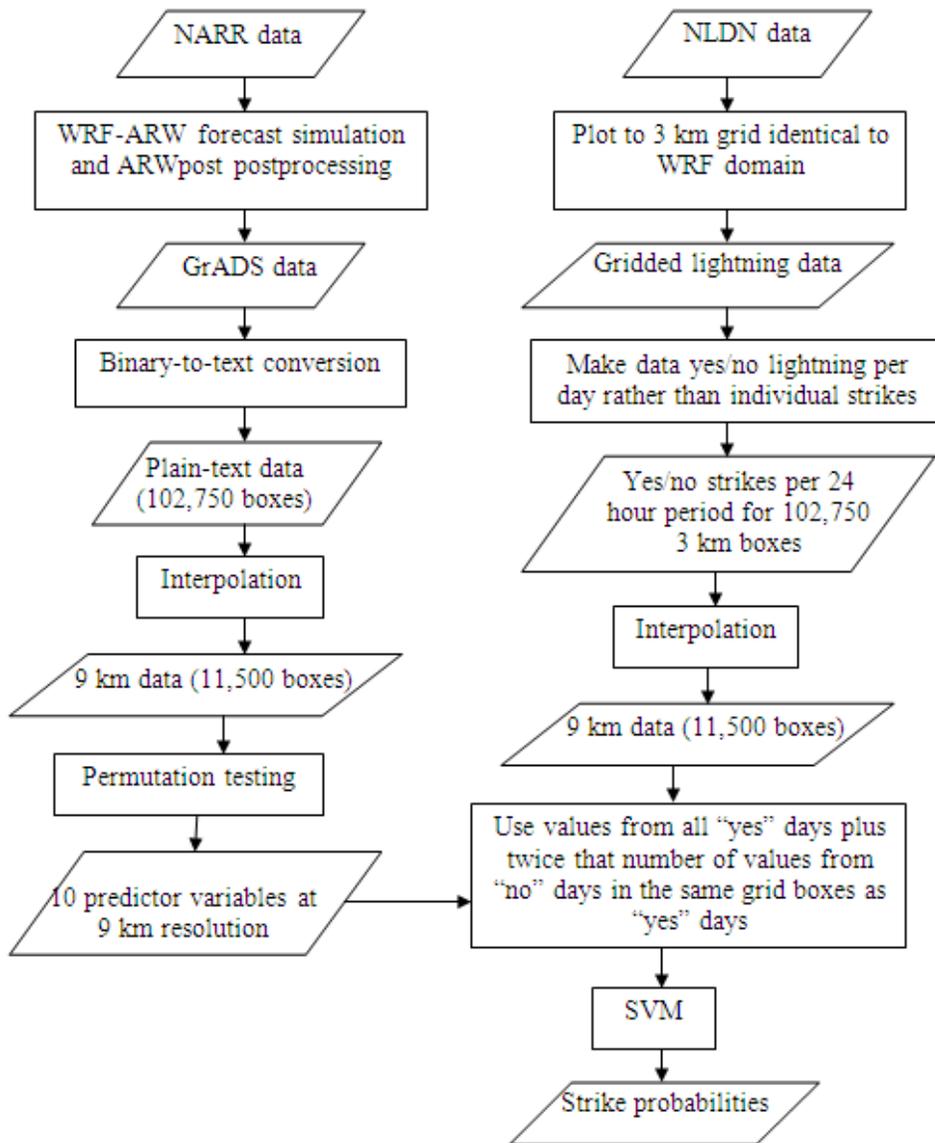
2.2.7 Validation of the SVM

Probabilistic lightning forecasts generated by the SVM were verified on a grid-scale level. Any lightning strike observations in a 9 km grid box are counted as a “yes,” and the absence of strikes in a grid box a “no.” These data points were compared against the SVM predictions for each grid box, and contingency statistics were formulated for a given day’s forecast. The predictor variables were cross-validated against randomly

selected individual cases from the 90 positive events. The predictions of the SVM for each grid box were compared against the observed lightning for that grid box, and contingency statistics were used to determine the skill of the SVM.

2.2.8 Methods summary

This chart summarizes the steps taken in conducting the research.



CHAPTER III

RESULTS

3.1 Prediction of a large data set

The study determined that the SVM had predictive power for a large data set consisting of the full table of predictor values itself. This data set was composed of 10 parameters for 52,158 grid boxes, or all yes events and twice the number of yes events as no events sampled from the same grid boxes as the yes events. Contingency statistics for the full predictor set are given in Table 3.1. The statistics in this table are based on a lightning probability threshold for a yes greater than 0.5.

Table 3.1 Contingency statistics for full predictor set

Percent correct (PC)	0.868
Critical success index (CSI)	0.653
Bias	0.877
False alarm ratio	0.155
Probability of detection	0.741
Probability of false detection	0.068
Heidke Skill Score	0.695
True skill statistic	0.674

As the table demonstrates, given a large data set for which to make predictions, the SVM was able to predict lightning with reasonably high accuracy and with high skill. The model does have a bias in favor of under-prediction of lightning.

3.2 Prediction of individual cases

While the model does reasonably well at predicting on the data from which it was trained, the model does not have skill with individual cases. A variety of these cases have been selected to illustrate this result. Unlike the training data set described in §3.1, which contains only those grid boxes in which lightning was recorded, these data sets include all 11,500 9 km grid boxes in the state of Mississippi. Therefore, the data sets for individual days will have a much lower percent coverage of lightning strikes than the training data set, which had 33 percent coverage. The day with the highest strike count of the 90 days used in this study, for example, had only 12.3 percent coverage.

3.2.1 3 October 2002

On 3 October 2002, there were 591 lightning strikes reported in Mississippi. A map of these strikes is shown in Figure 3.1. A map of the SVM's predicted lightning probability is shown in Figure 3.2. Contingency statistics for this case are shown in Table 3.2. The calculations in this table are based on a predicted probability of lightning greater than 0.5. As the figures and tables illustrate, the SVM did not show skill at predicting lightning on this day. This model vastly under-predicted lightning for this event, missing the eastern location of the lightning that occurred and under-predicting it in the area where it was placed instead.

Strikes for 3 October 2002

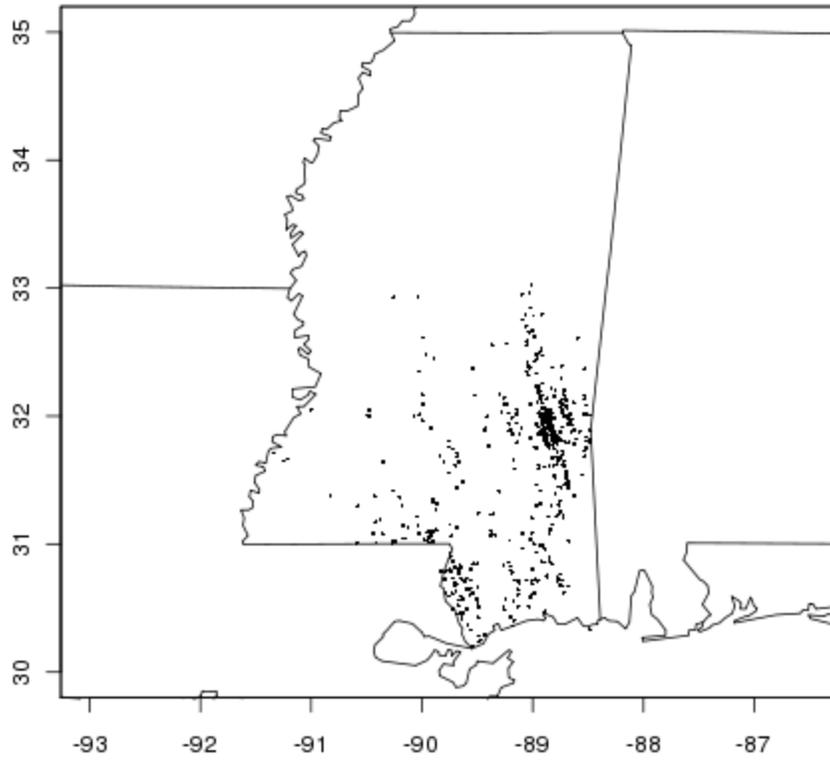


Figure 3.1 Observed lightning on 3 October 2002

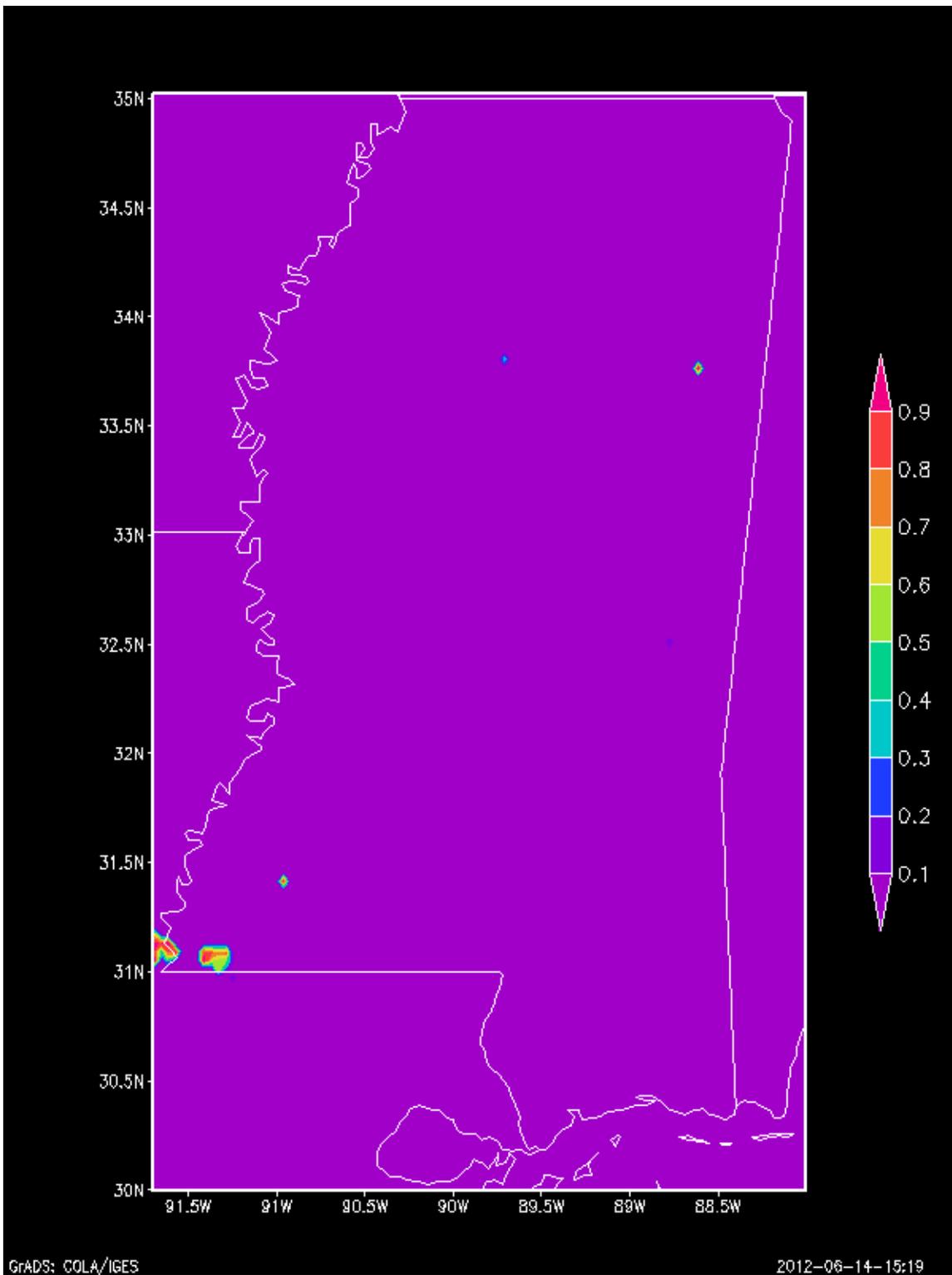


Figure 3.2 SVM-predicted lightning on 3 October 2002

Table 3.2 Contingency statistics for 3 October 2002

Percent correct (PC)	0.983
Critical success index (CSI)	0
Bias	0.101
False alarm ratio	1
Probability of detection	0
Probability of false detection	0.001
Heidke Skill Score	-0.0029
True skill statistic	-0.0016

3.2.2 24 November 2001

On 24 November 2001, there were 30,285 lightning strikes reported in Mississippi. This was the largest number of strikes reported in any of the 90 lightning cases selected for this research. A map of observed lightning for this date is shown in Figure 3.3. A map of the SVM's predicted lightning for this event is shown in Figure 3.4. Contingency statistics for the case are shown in Table 3.3. The calculations in this table are based on a predicted probability of lightning greater than 0.5. With the predictor variables that are used, the SVM completely failed to predict the lightning event that occurred this day.

Strikes for 24 November 2001

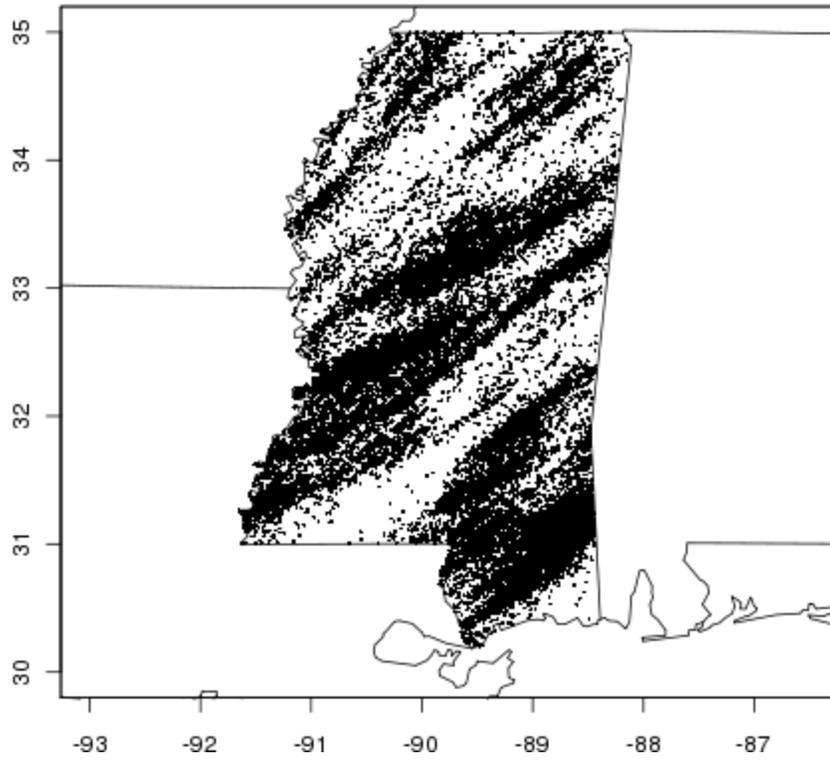


Figure 3.3 Observed lightning on 24 November 2001

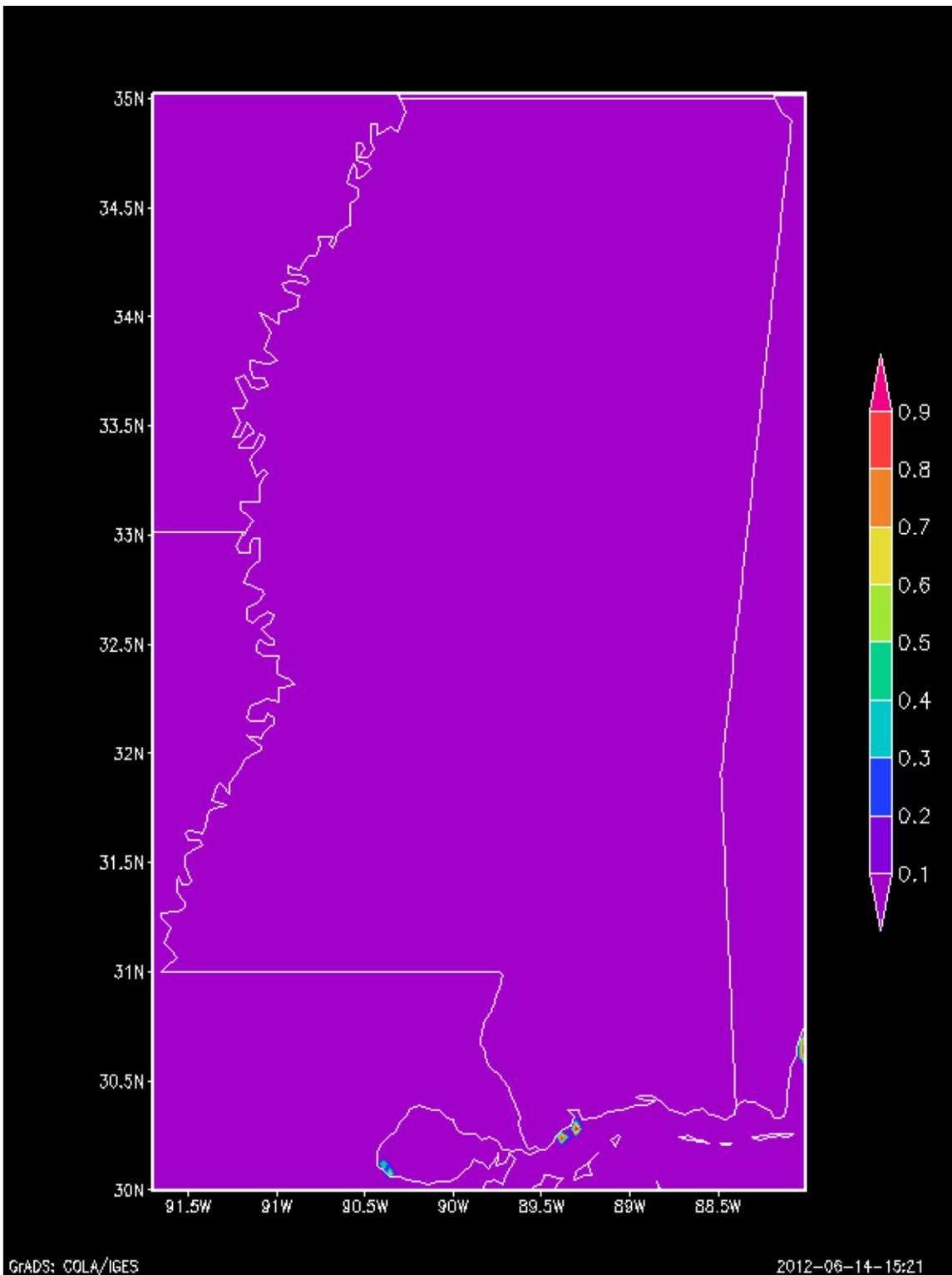


Figure 3.4 SVM-predicted lightning on 24 November 2001

Table 3.3 Contingency statistics for 24 November 2001

Percent correct (PC)	0.876
Critical success index (CSI)	0
Bias	0.0049
False alarm ratio	1
Probability of detection	0
Probability of false detection	0.0007
Heidke Skill Score	-0.0012
True skill statistic	-0.0007

3.2.3 20 October 2004

On 20 October 2004, there were 1,965 lightning strikes reported in Mississippi. A map of observed lightning for this date is shown in Figure 3.5. A map of the SVM's predicted lightning for this event is shown in Figure 3.6. Contingency statistics for the case are shown in Table 3.4. The calculations in this table are based on a predicted probability of lightning greater than 0.5. The SVM showed a percentage correct (PC) closer to the PC value for the full data set of 90 cases, but it still had a high probability of false detection and false alarm ratio. The SVM focused more lightning activity in the western area of Mississippi than the east, when actual lightning struck primarily in the eastern side of the state.

Strikes for 20 October 2004

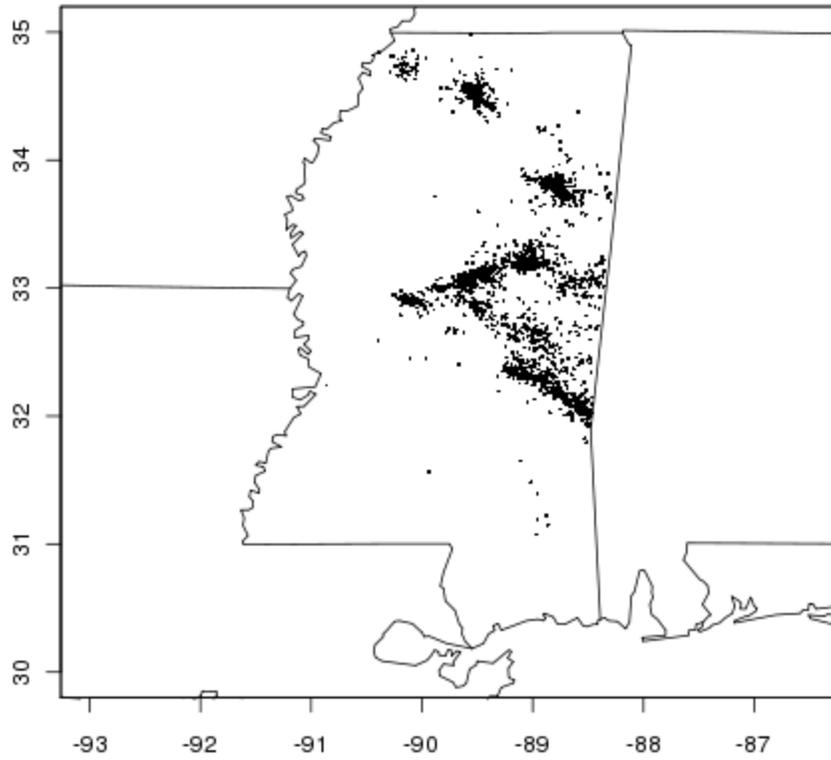


Figure 3.5 Observed lightning on 20 October 2004

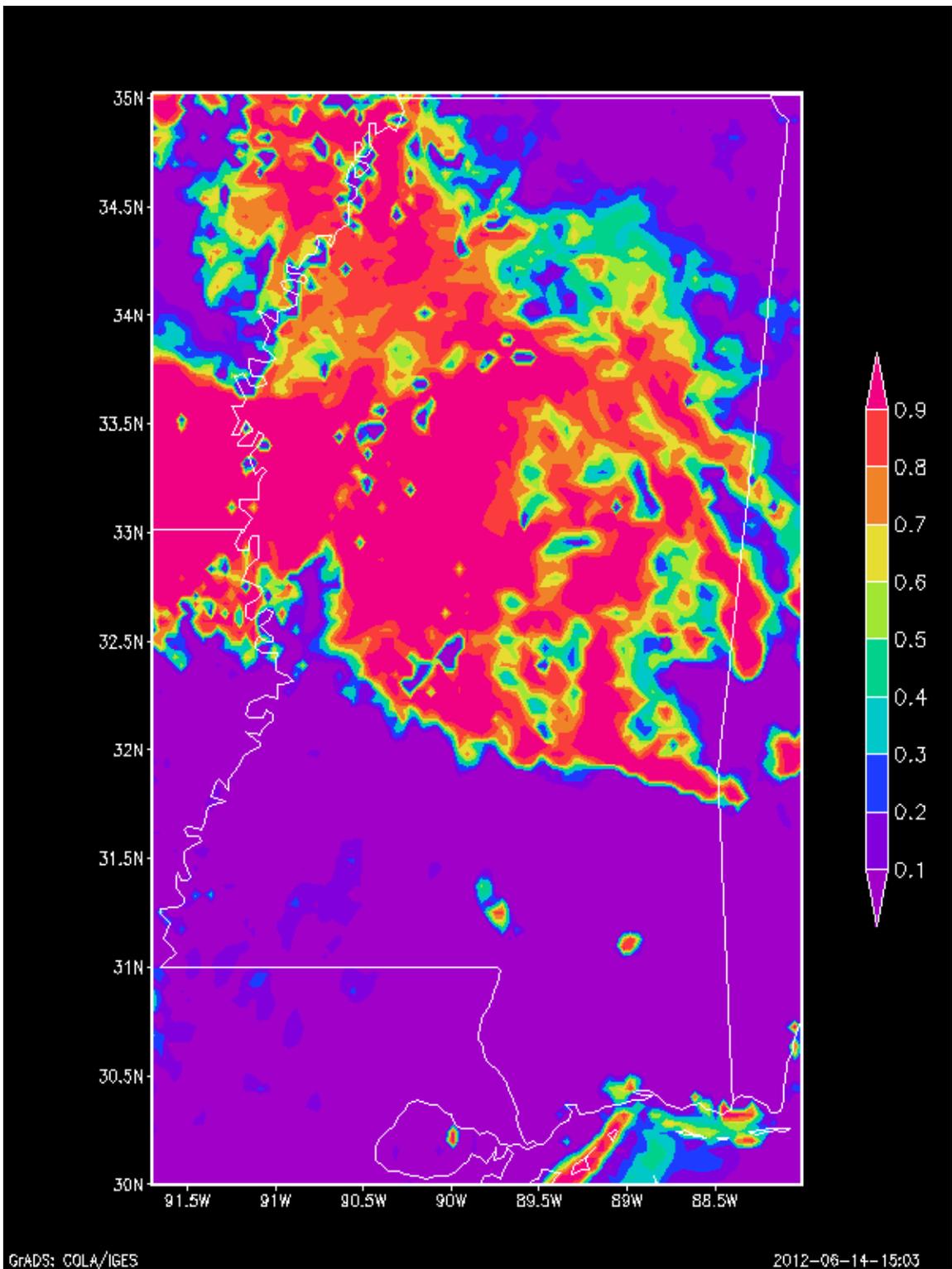


Figure 3.6 SVM-predicted lightning on 20 October 2004

Table 3.4 Contingency statistics for 20 October 2004

Percent correct (PC)	0.650
Critical success index (CSI)	0.052
Bias	13.3
False alarm ratio	0.947
Probability of detection	0.708
Probability of false detection	0.351
Heidke Skill Score	0.051
True skill statistic	0.357

3.2.4 22 September 2006

On 22 September 2006, there were 1,154 lightning strikes reported in Mississippi. A map of observed lightning for this date is shown in Figure 3.7. A map of the SVM's predicted lightning for this event is shown in Figure 3.8. Contingency statistics for the case are shown in Table 3.5. The calculations in this table are based on a predicted probability of lightning greater than 0.5. For this case, the same problem with under-prediction is present. The lightning areas on this date are suggestive of the tracks of discrete thunderstorms, the location of which would be inherently difficult to predict in advance.

Strikes for 22 September 2006

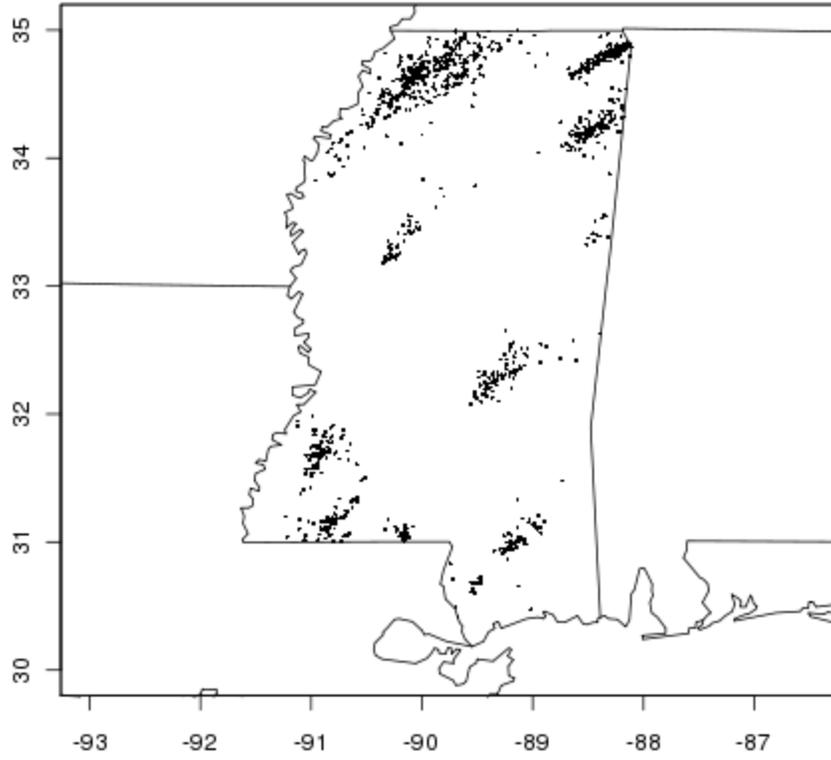


Figure 3.7 Observed lightning on 22 September 2006

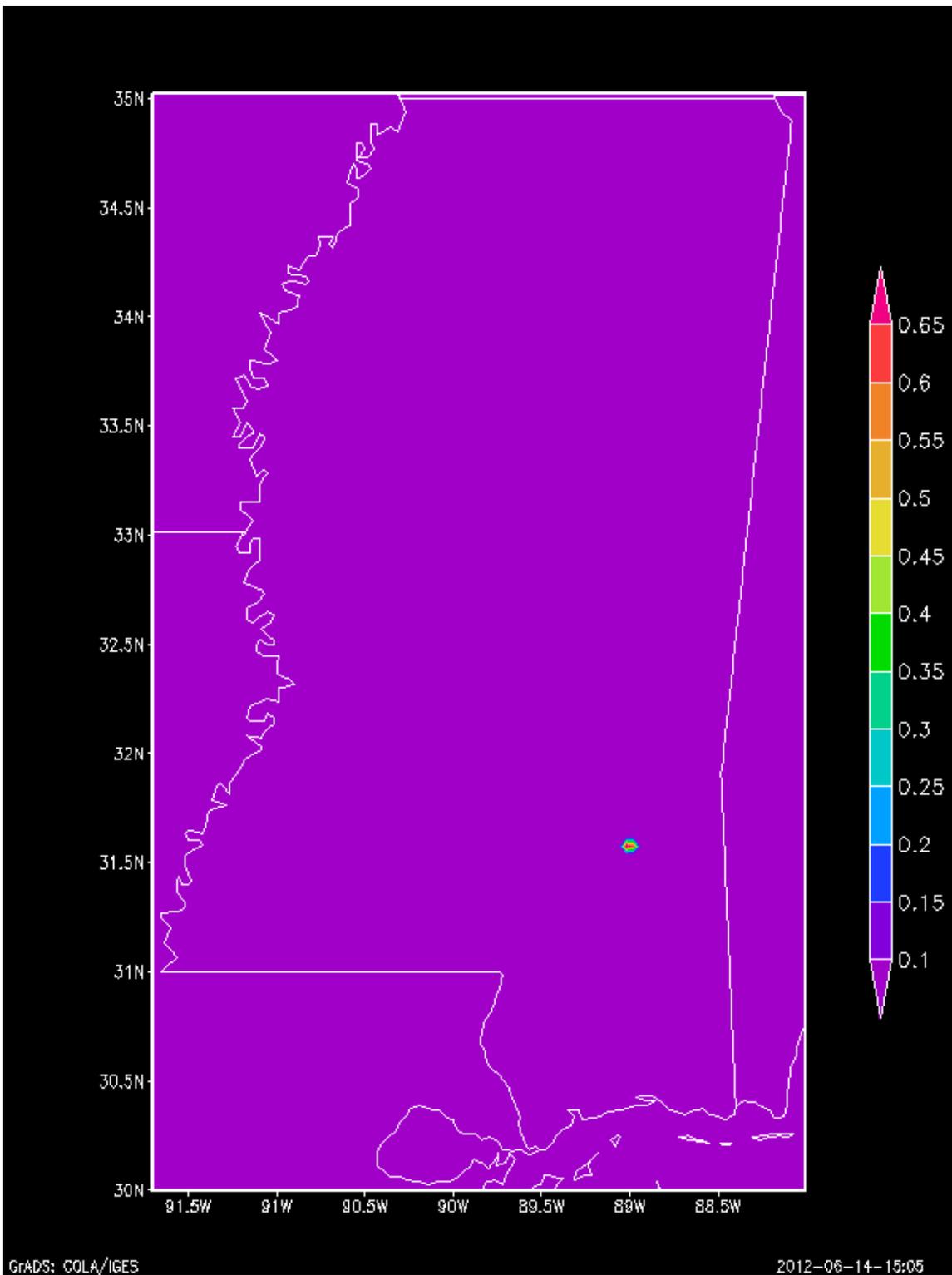


Figure 3.8 SVM-predicted lightning on 22 September 2006

Table 3.5 Contingency statistics for 22 September 2006

Percent correct (PC)	0.976
Critical success index (CSI)	0
Bias	0.0076
False alarm ratio	1
Probability of detection	0.000
Probability of false detection	0.0002
Heidke Skill Score	-0.0003
True skill statistic	-0.0002

3.2.5 16 September 2005

On 16 September 2005, there were 9,770 lightning strikes reported in Mississippi. A map of observed lightning for this date is shown in Figure 3.9. A map of the SVM's predicted lightning for this event is shown in Figure 3.10. Contingency statistics for the case are shown in Table 3.6. The calculations in this table are based on a predicted probability of lightning greater than 0.5. As the statistical table shows, this case too suffered from under-prediction of lightning. However, unlike the cases of 3 October 2002, 24 November 2001, and 22 September 2006, in which the SVM basically missed the event entirely, this SVM predicted lightning in areas that it did not occur. There also appears to be a spatial displacement of the weather system itself, as a comparison of observed lightning (Figure 3.9) and SVM-predicted lightning (Figure 3.10) suggests. Notably, an area of high lightning activity on this date in southwest Mississippi appears to be displaced into Louisiana in the SVM's predictions, and a secondary area of high lightning activity near Tupelo, Mississippi has been displaced to the far northeast corner in the SVM's predictions.

Strikes for 16 September 2005

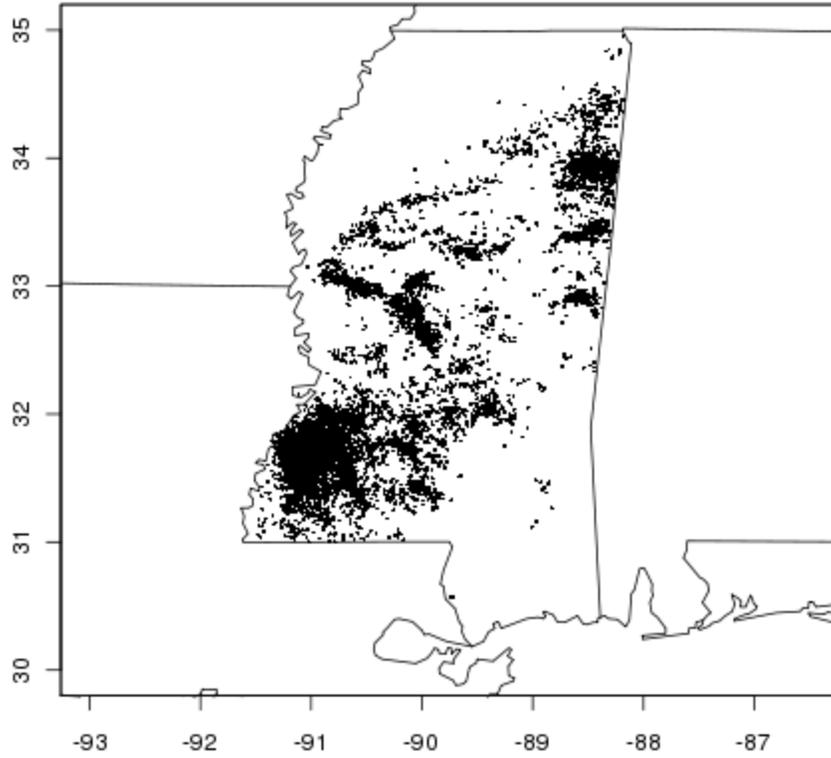


Figure 3.9 Observed lightning on 16 September 2005

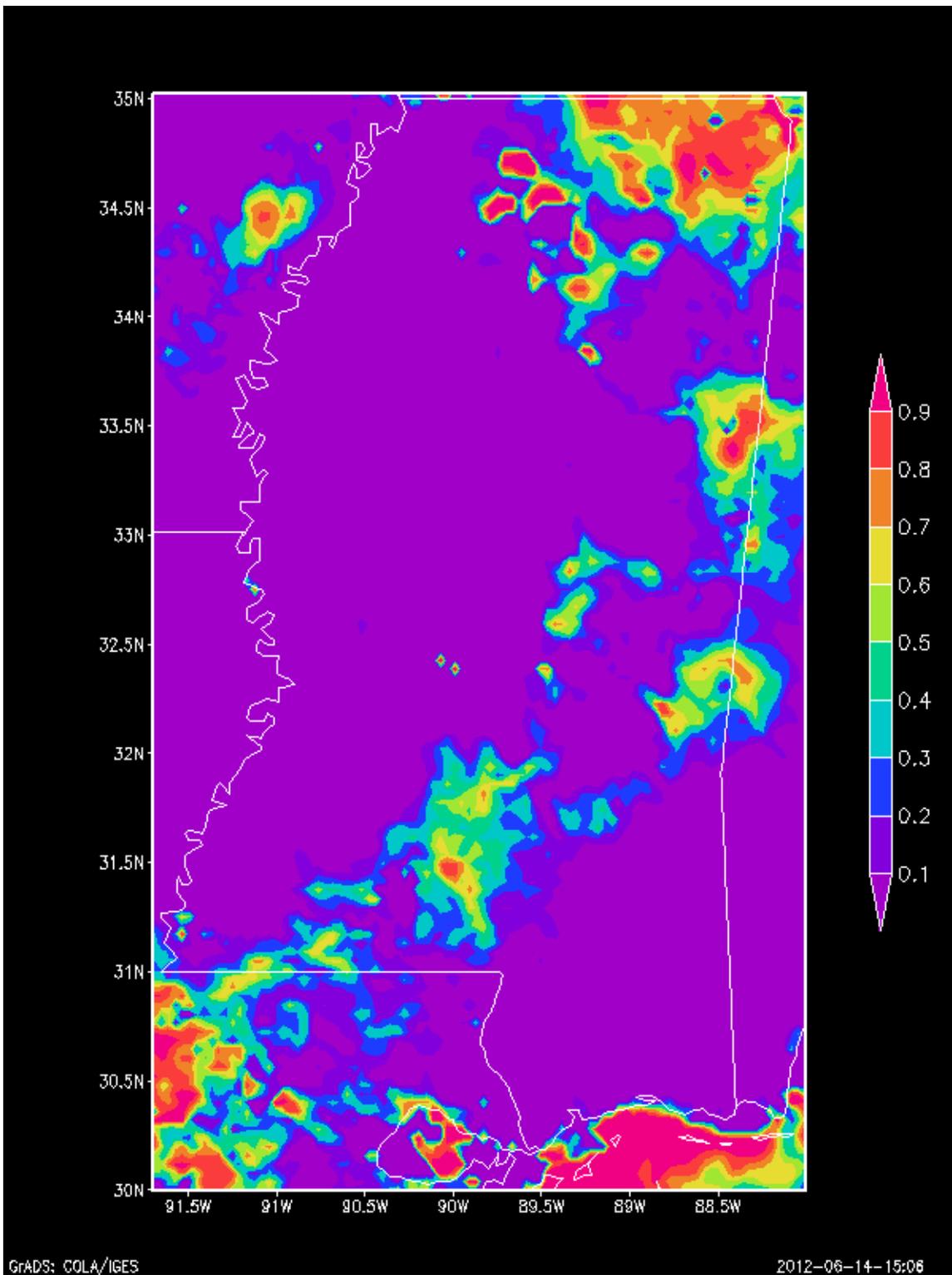


Figure 3.10 SVM-predicted lightning on 16 September 2005

Table 3.6 Contingency statistics for 16 September 2005

Percent correct (PC)	0.831
Critical success index (CSI)	0.020
Bias	2.154
False alarm ratio	0.971
Probability of detection	0.062
Probability of false detection	0.124
Heidke Skill Score	-0.040
True skill statistic	-0.061

3.3 Interpretation of results

The results are best described as mixed. The SVM showed predictive power for forecasting lightning over a data set consisting of 10 predictor values for 52,158 9 km grid boxes. However, it had limited predictive power for a smaller data set consisting of an individual 24-hour period. Many individual forecasts showed no skill at all. Overall the SVM under-predicted lightning for individual cases. Numerous factors may have contributed to the SVM's inability to forecast lightning with accuracy, and these potential factors are discussed further in this section.

3.3.1 NARR and WRF-ARW limitations

Two likely contributors to the SVM's case-wise failure to forecast lightning were the NARR data and the WRF model. All predictor values were taken from WRF model output that had been interpolated from reanalyzed NARR data sets. The model forecast was sensitive to parameterization of the WRF model itself, and it is possible that the parameterizations of model physics used in this study were not optimal for interpolating thunderstorm-level processes to a mesoscale grid.

It is not likely that the WRF simulation used in this research resulted in a significant temporal displacement of weather systems. The WRF model produced a forecast of reanalyzed data that had 3-hour temporal resolution; any model errors in the timing of weather systems would have been corrected with the next set of NARR input data from 3 hours later, preventing a “butterfly effect” or a small error in timing going unchecked and leading to a vast difference between reanalyzed data and WRF-simulated forecasts.

However, the same cannot be said of spatial displacement of systems. The NARR data used as WRF input have 32 km spatial resolution, and the WRF model interpolated these data down to 3 km. Considering that the average thunderstorm would barely occupy a single grid box on the raw NARR data and cloud-to-ground lightning does not necessarily strike the ground directly underneath its parent storm (i.e., a lightning strike may have occurred in a grid box adjacent to one containing a storm), a distinct possibility exists that the WRF model may have failed to accurately interpolate thunderstorm-level processes from an input data set of 32 km resolution.

It should be noted that permutation testing did not find any statistically significant difference between lightning-day and non-lightning-day values of *any* meteorological parameter subjected to the testing. The parameters chosen for the SVM were those with the statistically significantly smallest p-values. The lack of statistical significance of parameters directly associated with convective cloud activity, such as accumulated graupel, accumulated rain, and forecast maximum reflectivity, may be an artifact of the insufficiently fine resolution of the NARR data and the inability of the WRF model to accurately interpolate these data to a smaller grid.

Other operational forecast models such as the Rapid Update Cycle (RUC) and High Resolution Rapid Refresh (HRRR), which integrate observations from sources such as aircraft and live radar, may have greater accuracy at interpolating to a finely gridded spatial resolution.

3.3.2 Rarity of the event

Another factor that may have contributed to the SVM's inability to consistently accurately forecast lightning over a single 24-hour period is the rarity of lightning itself. Lightning is an inherently very unusual meteorological event, and its rarity is amply demonstrated by the fact that, when the state of Mississippi is divided in 102,750 3 km grid boxes, even on a high-impact severe weather day (24 November 2001), only 3,697 of these grid boxes had lightning observed. Even the interpolation to 11,500 9 km grid boxes yielded only 1,415 grid boxes on this day, which had the highest total number of lightning strikes in Mississippi of all the days selected for this research.

The option of running the SVM on lower-resolution inputs, for example, the native resolution of the NARR data (32 km), might seem to be a way to mitigate the rare-event issue. With a larger grid, the percent of grid boxes containing lightning would naturally increase. However, running the SVM model on a larger grid also increases the likelihood that individual thunderstorms and storm-scale processes will not be resolved well in the data set. The lack of statistically significant difference between atmospheric parameters when lightning occurred and when it did not occur again becomes a concern with a lower-resolution data set. Further study into these differences will be considered in future work.

CHAPTER IV

SUMMARY AND CONCLUSIONS

The purpose of this research was to examine the feasibility of a statistical learning-algorithmic approach for forecasting lightning probabilistically in high resolution, up to 1 day in advance of an anticipated event. The research used a support vector machine (SVM), a statistical machine learning algorithm, which was trained with reanalyzed and WRF-simulated atmospheric data from 90 days in meteorological autumn in which lightning occurred in the state of Mississippi. Permutation testing of lightning days against non-lightning days was used to determine the choice of predictor variables for the SVM. The SVM was trained with several possible kernels and cost functions, and little difference was noted between the performance of the different configurations of the model.

The research determined that an SVM does not, at present, have sufficient skill at forecasting lightning on a day-to-day basis to be used as an operational product. The model did show skill at forecasting lightning over a very large data set consisting of all 90 lightning days used in the study, with a Heidke Skill Score of 0.695. However, it showed much less skill, often no skill at all, at forecasting lightning on a statewide basis for a single forecast day.

Numerous factors may have contributed to the SVM's poor performance on a case-wise basis. The limitations of the WRF model itself in accurately interpolating 32

km NARR data to a 3 km grid likely generated some inaccuracy in the spatial distribution of thunderstorms and other mesoscale features. This possibility is supported by the fact that none of the atmospheric parameters subjected to permutation testing had statistically significant differences between lightning events and non-lightning events, even in a sample of the data specifically chosen to maximize the percentage of lightning events tested.

Another possibility for the poor performance of the SVM is the rarity of the event itself. Even high-impact thunderstorm days resulted in only fractional areas of the state (when subdivided into 3 km or 9 km grids) being impacted by lightning strikes. Interpolating the predictor values to a coarser grid would give the SVM a larger fraction of lightning events to work with, making the event appear less rare to the SVM, but this course of action would result in further loss of distinction between the values of atmospheric parameters when lightning is occurring and when it is not.

Despite the failure of the research to show that a support vector machine approach for forecasting CG lightning is ready for operational forecasting usage, the technique may still have potential for future research. Specifically, training SVMs with predictor data from other sources than a WRF-ARW simulation may have merit. Numerical forecast models with a higher spatial and temporal resolution, especially models that incorporate current conditions from unconventional but high-coverage sources such as radar reflectivity and aircraft measurements, may generate data that would provide an SVM with predictive power.

Such future work would also need to test different cost parameters to adjust the sensitivity of the SVM, different predictor variables, and different sampling techniques

for creating the training data set. With respect to predictor variables, the research found that none of the predictors generated by the WRF-ARW were statistically significantly different between yes and no events, but other numerical weather models might be able to produce predictors that had statistical significance (e.g., predictors directly related to cloud processes, such as graupel, rain, and simulated radar reflectivity, which such rapidly-updating forecast models would presumably model better than the WRF).

With respect to sampling of the training set, the training set used in this research contained 33 percent lightning events, which nonetheless resulted in an SVM that generally under-predicted lightning for actual weather events (with a percent coverage of 13 percent or lower). Finding an optimal percentage of lightning events for the training set would be a good research problem to address in the future.

REFERENCES

- Ackerman, T. P., A. J. Braverman, D. J. Diner, T. L. Anderson, R. A. Kahn, J. V. Martonchik, J. E. Penner, P. J. Rasch, B. A. Wielicki, and B. Yu, 2004: Integrating and interpreting aerosol observations and models within the PARAGON framework. *Bulletin of the American Meteorological Society*, Vol. 85, Issue 10, pp. 1523-1533.
- Aviolat, F., T. Cornu, and D. Cattani, 1998: Automatic clouds observation improved by an artificial neural network. *Journal of Atmospheric and Oceanic Technology*, Vol. 15, pp. 114-126.
- Berdeklis, P., and R. List, 2001: The ice-crystal-graupel collision charging mechanism of thunderstorm electrification. *Journal of the Atmospheric Sciences*, Vol. 58, pp. 2751-2770.
- Bright, D. R., M. S. Wandishin, R. E. Jewell, and S. J. Weiss, 2005: A physically based parameter for lightning prediction and its calibration in ensemble forecasts. *Preprints*, Conf. on Meteor. Applications of Lightning Data, San Diego CA.
- Burges, C. J. C., 1998: A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery* 2, pp. 121-167.
- Burrows, W. R., C. Price, and L. J. Wilson, 2005: Warm season lightning probability prediction for Canada and the northern United States. *Weather and Forecasting*, Vol. 20, pp. 971-988.
- Cummins, K. L., M. J. Murphy, E. A. Bardo, W. L. Hiscox, R. B. Pyle, and A. E. Pifer, 1998: A combined TOA/MDF technology upgrade for the U. S. National Lightning Detection Network. *Journal of Geophysical Research Atmospheres*, Vol. 103, No. D8, pp. 9035-9044.
- Cristianini, N., and J. Shawe-Taylor, 2000: *An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods*. Cambridge University Press, 189 pp.
- Dudhia, J., 1989: Numerical study of convection observed during winter monsoon experiment using a mesoscale two-dimensional model. *Journal of Atmospheric Science*, Vol. 46, pp. 3077-3107.

- , 1996: A multi-layer soil temperature model for MM5. *Preprints*, The Sixth PSU/NCAR Mesoscale Model Users' Workshop, 22-24 July 1996, Boulder, CO, pp. 49-50.
- Ebisuzaki, W., and G. Rutledge (ed.), 2004: Data documentation for NOAA Operational Model Archive and Distribution System (NOMADS) North America Regional Reanalysis (NARR) "Merge" data set. National Climatic Data Center.
- Efron, B., and R. Tibshirani, 1993: *An Introduction to the Bootstrap*. CRC Press, 436 pp.
- Fierro, A. O., M. S. Gilmore, E. R. Mansell, L. J. Wicker, and J. M. Straka, 2006: Electrification and lightning in an idealized boundary-crossing supercell simulation of 2 June 1995. *Monthly Weather Review*, Vol. 134, pp. 3149-3172.
- Fierro, A. O., L. M. Leslie, E. R. Mansell, and J. M. Straka, 2008: Numerical simulations of the microphysics and electrification of the weakly electrified 9 February 1993 TOGA COARE squall line: Comparisons with observations. *Monthly Weather Review*, Vol. 136, pp. 364-379.
- Hall, W. D., R. M. Rasmussen, and G. Thompson, 2005: The new Thompson microphysics scheme in WRF. WRF/MM5 Users' Workshop.
- Hallett, J., and C. P. R. Saunders, 1979: Charge separation associated with secondary ice crystal production. *Journal of the Atmospheric Sciences*, Vol. 36, pp. 2230-2235.
- Han, G., and Y. Shi, 2008: Development of an Atlantic Canadian coastal water level neural network model. *Journal of Atmospheric and Oceanic Technology*, Vol. 25, pp. 2117-2132.
- Hearst, M. A., S. T. Dumais, E. Osuna, J. Platt, and B. Schölkopf, 2002: Support vector machines. *Intelligent Systems and Their Applications, IEEE*, Vol. 13, Issue 4, pp. 18-28.
- Hong, S.-Y., and H.-L. Pan, 1996: Nonlocal boundary layer vertical diffusion in a medium-range forecast model. *Monthly Weather Review*, Vol. 124, pp. 2322-2339.
- Idone, V. P., D. A. Davis, P. K. Moore, Y. Wang, R. W. Henderson, M. Ries, and P. F. Jamason, 1998: Performance evaluation of the U. S. National Lightning Detection Network in eastern New York. Part II: Location accuracy. *Journal of Geophysical Research Atmospheres*, Vol. 103, No. D8, pp. 9057-9069.
- Jayarathne, E. R., and D. J. Griggs, 1991: Electric charge separation during the fragmentation of rime in an airflow. *Journal of the Atmospheric Sciences*, Vol. 48, No. 23, pp. 2492-2495.

- Jayarathne, E. R., C. P. R. Saunders, and J. Hallett, 1983: Laboratory studies of the charging of soft-hail during ice crystal interactions. *Quarterly Journal of the Royal Meteorological Society*, Vol. 109, pp. 609-630.
- Jennings, S. G., 1975: Charge separation due to water droplet and cloud droplet interactions in an electric field. *Quarterly Journal of the Royal Meteorological Society*, Vol. 101, pp. 227-234.
- Kitzmilller, D., M. A. R. Lilly, and S. D. Vibert, 2000: The SCAN 0-3 hour rainfall and lightning forecast algorithms. Office of Systems Development, National Weather Service.
- Koch, S. E., B. Ferrier, M. T. Stoelinga, E. Szoke, S. J. Weiss, and J. S. Kain, 2005: The use of simulated radar reflectivity fields in diagnosis of mesoscale phenomena from High-Resolution WRF model forecasts.” Preprints, *11th Conference on Mesoscale Processes and 32nd Conference on Radar Meteorology*, Albuquerque, NM.
- Lee, Y., G. Wahba, and S. A. Ackerman, 2004: Cloud classification of satellite radiance data by multcategory support vector machines. *Journal of Atmospheric and Oceanic Technology*, Vol. 21, pp. 159-169.
- Liu, H., V. Chandrasekar, and G. Xu, 2001: An adaptive neural network scheme for radar rainfall estimation from WSR-88D observations. *Journal of Applied Meteorology*, Vol. 40, pp. 2038-2050.
- McCaul Jr., E. W., S. J. Goodman, K. M. LaCasse, and D. J. Cecil, 2009: Forecasting lightning threat using cloud-resolving model simulations. *Weather and Forecasting*, Vol. 24, pp. 709-729.
- Marzban, C., and G. J. Stumpf, 1996: A neural network for tornado prediction based on Doppler radar-derived attributes. *Journal of Applied Meteorology*, Vol. 35, pp. 617-626.
- Mazany, R. A., S. Businger, S. I. Gutman, and W. Roeder, 2002: A lightning prediction index that utilizes GPS integrated precipitable water vapor. *Weather and Forecasting*, Vol. 17, pp. 1034-1047.
- Mercer, A. E., C. M. Shafer, C. A. Doswell III, L. M. Leslie, and M. B. Richman, 2009: Objective classification of tornadic and nontornadic severe weather outbreaks. *Monthly Weather Review*, Vol. 137, pp. 4355-4368.
- Mercer, A. E., M. B. Richman, and H. B. Bluestein, 2008: Statistical modeling of downslope windstorms in Boulder, Colorado. *Weather and Forecasting*, Vol. 23, pp. 1176-1194.

- Mesinger, F., G. DiMego, E. Kalnay, K. Mitchell, P. C. Shafran, W. Ebisuzaki, D. Jović, J. Woollen, E. Rogers, E. H. Berbery, M. B. Ek, Y. Fan, R. Grumbine, W. Higgins, H. Li, Y. Lin, G. Manikin, D. Parrish, and W. Shi, 2006: North American regional reanalysis. *Bulletin of the American Meteorological Society*, Vol. 87, pp. 343-360.
- Miller Jr., S. D., G. W. Carbin, J. S. Kain, E. W. McCaul, A. R. Dean, C. J. Melick, and S. J. Weiss, 2010: Preliminary investigation into lightning hazard prediction from high resolution model output. *Preprints*, 25th Conf. Severe Local Storms, Denver CO.
- Mlawer, E. J., S. J. Taubman, P. D. Brown, M. J. Iacono, and S. A. Clough, 1997: Radiative transfer for inhomogeneous atmosphere: RRTM, a validated correlated-k model for the longwave. *Journal of Geophysical Research Atmospheres*, V. 102 (D14), pp. 16663-16682.
- National Climatic Data Center (NCDC), 2003: Contents of NARR output AWIPS GRIB files.
- Orville, R. E., G. R. Huffines, W. R. Burrows, R. L. Holle, and K. L. Cummins, 2002: The North American Lightning Detection Network (NALDN)—First results: 1998-2000. *Monthly Weather Review*, Vol. 130, Issue 8, pp. 2098-2109.
- , 2008: Development of the National Lightning Detection Network. *Bulletin of the American Meteorological Society*, pp. 180-190.
- Reynolds, S. E., M. Brook, and M. F. Gourley, 1957: Thunderstorm charge separation. *Journal of Meteorology*, Vol. 14, pp. 426-436.
- Saunders, C. P. R., 1993: A review of thunderstorm electrification processes. *Journal of Applied Meteorology*, Vol. 32, pp. 642-655.
- Shih, C., 2010: The North American Regional Reanalysis (NARR) archive at NCAR. National Center for Atmospheric Research.
- Skamarock, W. C., and J. P. Klemp, 2007: A time-split nonhydrostatic atmospheric model for weather research and forecasting applications. *Journal of Computational Physics*, pp. 3465-3485.
- , ----, J. Dudhia, D. O. Gill, D. M. Barker, W. Wang, and J. G. Powers, 2005: A description of the Advanced Research WRF version 2. NCAR Tech. Note NCAR/TN-468+STR, 88 pp.
- Stackpole, J. D., 1994: The storage of weather product information and the exchange of weather product messages in gridded binary form. World Meteorological Organization, Pub. 306.

- Takahashi, T., 1978: Riming electrification as a charge generation mechanism in thunderstorms. *Journal of the Atmospheric Sciences*, Vol. 35, pp. 1536-1548.
- Thompson, Gregory, P. R. Field, W. D. Hall, and R. M. Rasmussen, 2006: A new bulk microphysical parameterization for WRF (& MM5). National Center for Atmospheric Research.
- Wilks, D. S., 2006: *Statistical Methods in the Atmospheric Sciences*. 2nd edition. Academic Press, 627 pp.
- Williams, E. R. R. Zhang, and J. Rydock, 1991: Mixed-phase microphysics and cloud electrification. *Journal of the Atmospheric Sciences*, Vol. 48, No. 19, pp. 2195-2203.

APPENDIX A
CASES USED IN THE RESEARCH

Table A.1 lists the cases chosen for the research and the number of lightning strikes reported on each day (0000 UTC to 2359 UTC) in Mississippi.

Table A.1 Lightning cases and total lightning strikes in Mississippi

Case	Total Strikes		
2001-09-02	2682	2003-09-07	932
2001-09-03	3715	2003-09-09	1271
2001-09-04	5630	2003-09-13	3117
2001-09-05	1498	2003-10-06	15
2001-09-13	1084	2003-10-10	53
2001-09-17	32	2003-11-12	49
2001-09-18	72	2003-11-24	44
2001-09-19	1938	2004-09-02	357
2001-09-21	11	2004-09-11	16
2001-10-10	162	2004-09-14	12
2001-10-14	28	2004-10-10	379
2001-10-25	5394	2004-10-12	216
2001-10-26	1	2004-10-19	20659
2001-11-14	1	2004-10-20	1965
2001-11-24	30285	2004-10-25	746
2001-11-25	27	2004-10-28	222
2001-11-27	10695	2004-10-29	1
2001-11-29	22293	2004-11-02	1259
2002-09-01	1	2004-11-04	10
2002-09-03	51	2004-11-20	23
2002-09-07	1398	2004-11-21	598
2002-09-09	3	2004-11-22	746
2002-09-16	1193	2004-11-23	1746
2002-10-01	89	2005-09-01	4
2002-10-03	591	2005-09-04	16
2002-10-04	19	2005-09-15	373
2002-10-05	1215	2005-09-16	9770
2002-10-07	12047	2005-09-17	6
2002-10-09	225	2005-09-29	1068
2002-10-10	691	2005-10-01	1503
2002-10-15	192	2005-10-08	1
2002-10-28	6618	2005-10-21	1
2002-11-15	1269	2006-09-07	46
2002-11-26	17	2006-09-14	4
		2006-09-18	2804

Table A.1 (continued)

2006-09-19	557
2006-09-22	1154
2006-10-02	328
2006-10-17	3788
2006-10-18	228
2006-10-21	62
2006-10-27	168
2006-11-01	5593
2006-11-06	830
2007-09-04	3172
2007-09-05	756
2007-09-06	89
2007-09-09	217
2007-09-10	99
2007-10-01	1
2007-10-04	8
2007-10-15	17
2007-11-06	699
2007-11-14	969
2007-11-15	1607
2007-11-22	6998