# Project Report

### Pengfei Gao and Tim Moon

June 8, 2016

## 1 Date Processing

We investigate several machine learning models to predict temperature based on date, location, and weather conditions. These models are trained with a data set from the National Oceanic and Atmospheric Administration. This data set ranges from 1929 to 2016 and consists of 139 million daily weather reports from weather stations throughout the world. We are interested in predicting average temperature based on the following features:

Feature	Type	Units	Comments
Station ID	Categorical		Concatenation of DATSAV3 ID and WBAN ID
Year	Numeric		
Month	Numeric		
Day	Numeric		
Dew Point	Numeric	$^{\circ}\mathrm{F}$	
Sea Level Pressure	Numeric	$_{ m mbar}$	
Station Pressure	Numeric	$_{ m mbar}$	
Visibility	Numeric	$_{ m mi}$	
Mean Wind Speed	Numeric	${\rm kn}$	
Maximum Wind Seed	Numeric	kn	
Gust Speed	Numeric	${ m kn}$	
Precipitation	Numeric	in	
Precipitation Report	Categorical		Indicates procedure to measure precipitation
Fog	Categorical		Boolean
Rain	Categorical		Boolean
$\operatorname{Snow}$	Categorical		Boolean
Hail	Categorical		Boolean
Thunder	Categorical		Boolean
Tornado	Categorical		Boolean

# 2 Implementation

The data set was randomly split into training (70%) and validation (30%) sets and the training set was used to train ordinary least squares, generalized linear, gradient boosting, and random forest estimators. We also tried a grid search for hyper-parameters of gradient boosting algorithm. This was implemented with H2O to allow for quick prototyping. In particular, H2O implements distributed data structures and algorithms, allowing us to easily leverage parallel computer architectures. Experiments were performed using four EC2 instances from Amazon Web Services.

#### 3 Results

The mean squared errors for each method are reported below:

Model	Train MSE	Train $\mathbb{R}^2$	Validation MSE	Validation $\mathbb{R}^2$
Ordinary Least Squares	52.7	0.91	52.7	0.91
Generalized Linear Model	99.7	0.83	99.6	0.83
Gradient Boosting	138.0	0.76	137.8	0.76
Random Forest	45.3	0.92	43.8	0.92

We see that the random forest method yields the best performance. We also note that ordinary least squares outperformed the generalized linear model, likely since it used 26716 predictors compared to the generalized linear model's 28. We hypothesize that the gradient boosting method struggled because the data set exhibits behavior that is too complicated to be captured by shallow trees. Although the performance can improve by incorporating more trees, it may take an excessive number to achieve a good model. To investigate this, we performed a hyperparameter search with the gradient boosting method:

Maximum Tree Depth	Learning Rate	Residual Deviance
5	0.1	137.8
5	0.05	138.1
3	0.1	146.4
3	0.05	146.6
5	0.01	299.0
3	0.01	310.6

We see that the performance tends to improve as the trees become deeper, consistent with our hypothesis. In addition, the performance increases as the learning rate becomes larger. Inspecting the random forest, we find that the most important features are as follows:

Feature	Scaled Importance
Dew Point	1
Station ID	0.16
Snow	0.12
Sea Level Pressure	0.09
Month	0.08

Dew point is a measure of humidity, which indicates that temperature is very sensitive to humidity. The standardized coefficients computed by ordinary least squares are as follows (excluding non-Boolean categorical features):

<sup>&</sup>lt;sup>1</sup>Recall that categorical features must be converted to multiple numerical features before applying a generalized linear model. Thus, it appears that the generalized linear model discarded most of the information involving station ID.

Year         0.576           Month         0.0760           Day         0.0025           Dew Point         20.5           Sea Level Pressure         -1.69           Station Pressure         0.0035           Visibility         1.377           Mean Wind Speed         0.244           Max Wind Speed         0.449           Gust Speed         -0.467           Precipitation         -0.436           Snow Depth         -0.424           Fog         -3.04           Rain         -3.14           Snow         -5.29           Hail         1.23           Thunder         -0.07           Tornado         -1.10	Feature	Standardized Coefficient	
Day       0.0025         Dew Point       20.5         Sea Level Pressure       -1.69         Station Pressure       0.0035         Visibility       1.377         Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Year	0.576	
Dew Point       20.5         Sea Level Pressure       -1.69         Station Pressure       0.0035         Visibility       1.377         Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Month	0.0760	
Sea Level Pressure       -1.69         Station Pressure       0.0035         Visibility       1.377         Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Day	0.0025	
Station Pressure       0.0035         Visibility       1.377         Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Dew Point	20.5	
Visibility       1.377         Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Sea Level Pressure	-1.69	
Mean Wind Speed       0.244         Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Station Pressure	0.0035	
Max Wind Speed       0.449         Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Visibility	1.377	
Gust Speed       -0.467         Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Mean Wind Speed	0.244	
Precipitation       -0.436         Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Max Wind Speed	0.449	
Snow Depth       -0.424         Fog       -3.04         Rain       -3.14         Snow       -5.29         Hail       1.23         Thunder       -0.07	Gust Speed	-0.467	
Fog -3.04 Rain -3.14 Snow -5.29 Hail 1.23 Thunder -0.07	Precipitation	-0.436	
Rain -3.14 Snow -5.29 Hail 1.23 Thunder -0.07	Snow Depth	-0.424	
Snow -5.29 Hail 1.23 Thunder -0.07	Fog	-3.04	
Hail 1.23 Thunder -0.07	Rain	-3.14	
Thunder -0.07	$\operatorname{Snow}$	-5.29	
	Hail	1.23	
Tornado -1.10	Thunder	-0.07	
	Tornado	-1.10	

Observe that the temperature is positively correlated with the year, which is suggestive of global warming.

## Appendix A Python Script

```
1 import csv
2 import h2o
3 import sys
4 orig_stdout = sys.stdout
6 # filepath = '/Users/Pengfei/Documents/data/'
7 filepath = 's3n://stanford-cme250a/weather/data/'
8 filelist =[ filepath+'Xheader.csv' ]
9 for year in range (1940, 2017):
         filelist.append(filepath+'X'+str(year)+'.csv')
10
12 # Initialize H20
13 # h2o.init()
14 h2o.init(ip="localhost",port=55555,strict_version_check=False)
16 # Load data from files
data = h2o.import_file(filelist)
20 # Convert Boolean data to categorical
21 data['fog'] = data['fog'].asfactor()
22 data['rain'] = data['rain'].asfactor()
23 data['snow'] = data['snow'].asfactor()
24 data['hail'] = data['hail'].asfactor()
25 data['thunder'] = data['thunder'].asfactor()
26 data['tornado'] = data['tornado'].asfactor()
28 # Delete unnecessary columns
29 data = data.drop('tempucnt')
```

```
30 data = data.drop('dewpoint cnt')
31 data = data.drop('sea cnt')
32 data = data.drop('stat cnt')
33 data = data.drop('visiucnt')
data = data.drop('winduspeeducnt')
data = data.drop('*is,hourly,max')
data = data.drop('*is,hourly,min')
_{38} # Combine WBAN and DAVSAT3 station IDs to get unified station IDs
39 StationIds = 100000*data['Station'] + data['WBAN']
40 StationIds = StationIds.asfactor()
StationIds.set_name(0, 'StationId')
data = data.cbind(StationIds)
43 data = data.drop('Station')
44 data = data.drop('WBAN')
46 # Obtain date
months = data['MonthDay'] // 100
48 days = data['MonthDay'] % 100
49 months.set_name(0, 'Month')
50 days.set_name(0, 'Day')
51 data = data.cbind(months)
52 data = data.cbind(days)
53 data = data.drop('MonthDay')
55 # Remove entries with missing temperature data
56 data[data['temp']>9999,'temp'] = None
57 # data = data.na_omit()
59 # Remove missing data
60 data[data['dewpoint']>9999, 'dewpoint'] = None
61 data[data['sea_level_pres']>9999,'sea_level_pres'] = None
data[data['stationupres']>9999, 'stationupres'] = None
data[data['visibility']>999,'visibility'] = None
64 data[data['mean_wind_speed'] > 999, 'mean_wind_speed'] = None
65 data[data['max_wind_speed']>999,'max_wind_speed'] = None
66 data[data['gust_speed']>999, 'gust_speed'] = None
67 data[data['max_temp']>9999, 'max_temp'] = None
68 data[data['min_temp']>9999, 'min_temp'] = None
69 data[data['precipitation']>99, 'precipitation'] = 0
70 data[data['snow,depth']>999, 'snow,depth'] = 0
74 # Generate train set and validation set
75 [train, val] = data.split_frame(ratios=[0.7])
77 # set chosen feature
78 feature_list = list(data.names)
79 feature_list.remove('temp')
80 feature_list.remove('max utemp')
81 feature_list.remove('min temp')
82
83
```

```
84 # Training Models
86 glm = h2o.estimators.glm.H2OGeneralizedLinearEstimator(model_id='glm1')
87 glm.train(y = "temp", x = feature_list, training_frame = train, validation_frame = val
ss orig_stdout = sys.stdout
s9 outputFile = open('log.txt', 'a')
90 sys.stdout = outputFile
92 print glm
93 outputFile.close()
94 sys.stdout = orig_stdout
96 gbm = h2o.estimators.gbm.H2OGradientBoostingEstimator(model_id='gbm1', distribution='g
97 gbm.train(y = "temp", x = feature_list, training_frame = train, validation_frame = val
98 orig_stdout = sys.stdout
outputFile = open('log.txt', 'a')
sys.stdout = outputFile
102 print gbm
outputFile.close()
sys.stdout = orig_stdout
106 # Try grid search
from h2o.grid.grid_search import H2OGridSearch
108 hyper_parameters = {'ntrees':[50], 'max_depth':[3,5], 'learn_rate':[0.01,0.1,0.05]}
109 gs = H2OGridSearch(h2o.estimators.gbm.H2OGradientBoostingEstimator(distribution='gauss
110 gs.train(y = "temp", x = feature_list, training_frame = train, validation_frame = val)
orig_stdout = sys.stdout
outputFile = open('log.txt', 'a')
sys.stdout = outputFile
115 print gs
outputFile.close()
sys.stdout = orig_stdout
119 rf = h2o.estimators.random_forest.H2ORandomForestEstimator(model_id='rf1')
120 rf.train(y = "temp", x = feature_list, training_frame = train, validation_frame = val)
orig_stdout = sys.stdout
outputFile = open('log.txt', 'a')
sys.stdout = outputFile
125 print rf
outputFile.close()
sys.stdout = orig_stdout
129 #deep learning is extremely slow, might not include it in the big data version
130 dl = h2o.estimators.deeplearning.H2ODeepLearningEstimator(model_id='dl1')
131 dl.train(y = "temp", x = feature_list, training_frame = train, validation_frame = val)
132 orig_stdout = sys.stdout
outputFile = open('log.txt', 'a')
sys.stdout = outputFile
136 print dl
outputFile.close()
```