



Purple Air: Air Quality Analysis

Eden Ehm
BIA 742: Predictive Analytics
Spring 2019

+ What is Air Quality?

- The degree to which the ambient air is pollution free, as assessed by measuring a number of indicators of pollution.
- Implications for:
 - Health
 - Respiratory System
 - Skin
 - Exercise
 - Lifestyle
 - Safety



+ Decorah, Iowa's Air Quality

- We have really bad air quality!
- **Why?** Hypotheses:
 - Geography
 - Truck Traffic
 - Concentrated Feeding Operations
 - Farms and Agriculture
 - City Design
 - Weather



+ Question



- **What affects air quality?**
- **Can we predict if air quality is safe using only data readily available on any basic weather report?**
(Temperature, Humidity, Date, Location, etc.)

+ Question – Why?

- We want to answer these questions for several reasons:
 - Better understanding of what does or does not affect air quality
 - What is in our control? What is out of our control?
 - Knowing this has Environmental and Ecological implications
 - Discovering if there are connections between the different sizes of particulate matter in the air
 - What is the pollution?
 - Better lifestyle and health choices
 - Not everyone has their own air quality sensor or ability to look up air quality online. Can we get a rough idea of air quality (safe vs. unsafe) based on basic weather information?



+ Terminology



- **Air Quality Index:** The Environmental Protection Agency (EPA) calculates the AQI for five major air pollutants regulated by the Clean Air Act: ground-level ozone, **particle pollution (also known as particulate matter)**, carbon monoxide, sulfur dioxide, and nitrogen dioxide. For each of these pollutants, EPA has established national air quality standards to protect public health.

<https://airnow.gov/index.cfm?action=aqibasics.aqi>

- **Particulate Matter (PM):** Particulate matter is the sum of all solid and liquid particles suspended in air many of which are hazardous. This complex mixture includes both organic and inorganic particles, such as dust, pollen, soot, smoke, and liquid droplets. These particles vary greatly in size, composition, and origin. They can be directly emitted or indirectly formed.

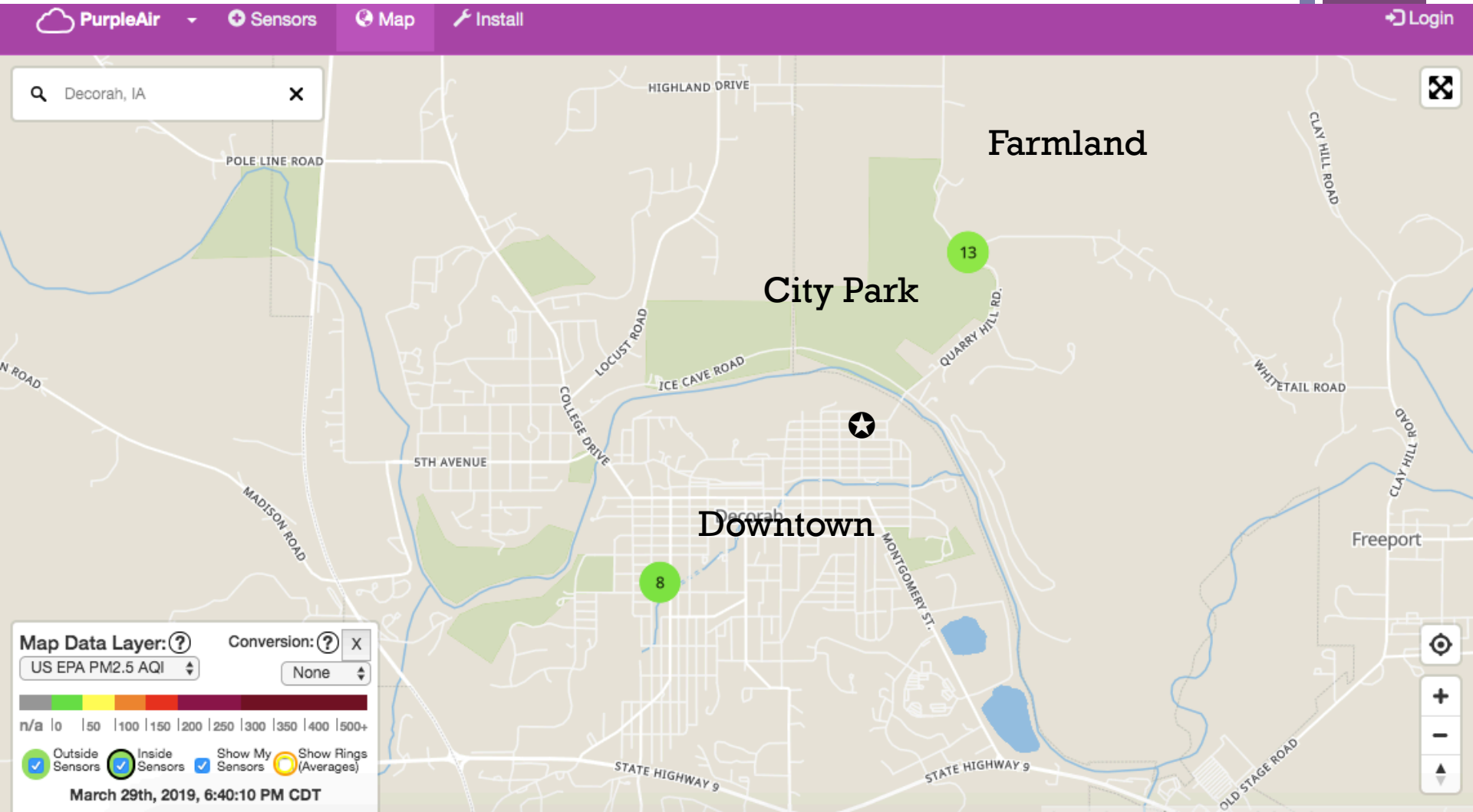
<https://www.greenfacts.org/en/particulate-matter-pm/level-2/01-presentation.htm>

+ Data Source

- Local Purple Air sensors
 - Decorah, IA – downtown
 - Decorah 2 – rural
- Utilize PMS5003 laser optical particle counters to collect air quality information from their surroundings. These sensors count suspended particles in sizes of 0.3, 0.5, 1.0, 2.5, 5.0 and 10 μ m. These particle counts are processed by the sensor using a complex algorithm to calculate the PM1.0, PM2.5 and PM10 mass in ug/m³.
- Also collect temperature, humidity, date, and other data



+ Data Source



+ Preparing the .csv Files



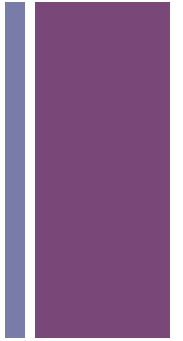
- Downloaded .csv files for the “Decorah 2” (rural) and “Decorah, IA” (downtown) sensors
- 1 Year: 4/1/2018 through 4/1/2019
- Combined into one .csv file called **PurpleAirDataCSV.csv**
 - Easy to do with Copy/Paste on my computer
- **PurpleAirDataCSV.csv** raw data file is attached to this submission

+ Preparing the .csv Files



- Added binary variable **location_id**
 - 1 = Decorah, IA downtown sensor
 - 0 = Decorah 2 rural sensor
- Added character variable **location**
 - “downtown” or “rural” classification
- Removed a duplicative column for PM2.5_CF_ATM_ug/m3
- Created **Total_PM** variable
 - Sum of PM1.0, PM2.5, and PM10.0
 - Not a “scientific” measurement in terms of air quality
 - A basic measure of “all of the particulate matter of all sizes in the air”

+ Adding a Binary Decision Variable



- Binary Variable called **Safe**
 - 1 = safe air quality
 - 0 = unsafe air quality (Total_PM >50)
- Purple Air sensors use US EPA PM2.5 AQI as their air quality index.
 - 0-50 safe
 - 50-100 acceptable
 - 100+ dangerous
- I chose 50 as the cutoff because it is *just* PM pollution.
It is better to be conservative with one's health and safety!
- 50 indicates that there is enough PM in the air that individuals with compromised respiratory systems will have issues & individuals with healthy systems might experience trouble while exercising.

+ Adding a Binary Decision Variable

Total_PM > 50 means think twice before going outside!



Air Quality Index (AQI) Values	Levels of Health Concern	Colors
<i>When the AQI is in this range:</i>	<i>..air quality conditions are:</i>	<i>...as symbolized by this color:</i>
0 to 50	Good	Green
51 to 100	Moderate	Yellow
101 to 150	Unhealthy for Sensitive Groups	Orange
151 to 200	Unhealthy	Red
201 to 300	Very Unhealthy	Purple
301 to 500	Hazardous	Maroon

+ Adding a Binary Decision Variable

```
*Add Binary DV;  
data ProjData.PurpleAirDataSAS;  
  set ProjData.PurpleAirDataSAS;  
  if Total_PM < 50 then Safe=1;  
  else Safe=0;  
  
run;
```

+ Exploring the Data

- proc contents
- 725,687 observations
- No missing data

	location_id	location	created_at_String	created_at	entry_id	PM1.0	PM2.5	PM10.0	UptimeMinutes	ADC	Temperature_F	Humidity	Total_PM	Safe
1	1	downtown	2018-04-01 00:00:56 UTC	01APR2018:00:00:56	527938	1.98	2.43	2.81	319	-75	35	39	7.22	0
2	1	downtown	2018-04-01 00:02:17 UTC	01APR2018:00:02:17	527939	3.21	4.35	4.37	320	-76	35	39	11.93	0
3	1	downtown	2018-04-01 00:03:36 UTC	01APR2018:00:03:36	527940	6.21	7.84	9.35	321	-79	35	40	23.40	0
4	1	downtown	2018-04-01 00:04:56 UTC	01APR2018:00:04:56	527941	4.29	5.07	5.14	323	-78	34	40	14.50	0
5	1	downtown	2018-04-01 00:06:16 UTC	01APR2018:00:06:16	527942	2.46	3.31	4.28	324	-77	33	40	10.05	0
6	1	downtown	2018-04-01 00:07:36 UTC	01APR2018:00:07:36	527943	1.74	3.16	4.09	325	-77	34	40	8.99	0
7	1	downtown	2018-04-01 00:08:56 UTC	01APR2018:00:08:56	527944	2.27	3.50	3.86	327	-75	33	40	9.63	0
8	1	downtown	2018-04-01 00:10:16 UTC	01APR2018:00:10:16	527945	2.38	3.21	3.21	328	-74	34	41	8.80	0
9	1	downtown	2018-04-01 00:11:36 UTC	01APR2018:00:11:36	527946	1.72	2.74	3.70	329	-75	34	41	8.16	0
10	1	downtown	2018-04-01 00:12:56 UTC	01APR2018:00:12:56	527947	1.93	3.05	3.30	331	-74	34	41	8.28	0
11	1	downtown	2018-04-01 00:14:16 UTC	01APR2018:00:14:16	527948	2.56	3.49	4.63	332	-76	34	41	10.68	0
12	1	downtown	2018-04-01 00:15:36 UTC	01APR2018:00:15:36	527949	2.44	2.98	3.49	333	-74	33	41	8.91	0
13	1	downtown	2018-04-01 00:16:57 UTC	01APR2018:00:16:57	527950	1.82	3.18	3.91	335	-76	33	41	8.91	0
14	1	downtown	2018-04-01 00:18:16 UTC	01APR2018:00:18:16	527951	1.88	3.22	3.17	336	-75	33	41	8.16	0

+ Exploring the Data

Alphabetic List of Variables and Attributes						
#	Variable	Type	Len	Format	Informat	Label
10	ADC	Num	8	F3.		ADC
12	Humidity	Num	8	F3.		Humidity
8	PM10_0_CF_ATM_ugm3	Num	8	F7.2		PM10.0_CF_ATM_ugm3
6	PM1_0_CF_ATM_ugm3	Num	8	F6.2		PM1.0_CF_ATM_ugm3
7	PM2_5_CF_ATM_ugm3	Num	8	F7.2		PM2.5_CF_ATM_ugm3
14	Safe	Num	8	F1.		Safe
11	Temperature_F	Num	8	F4.		Temperature_F
13	Total_PM	Num	8	F7.2		Total_PM
9	UptimeMinutes	Num	8	F5.		UptimeMinutes
4	created_at	Num	8	DATETIME19.		created_at
3	created_at_String	Char	23	\$23.	\$23.	created_at_String
5	entry_id	Num	8	F6.		entry_id
2	location	Char	8	\$8.	\$8.	location
1	location_id	Num	8	F1.		location_id

+ Exploring the Data: Looking for Problems

- `proc univariate`
- **PM1.0, PM2.5, PM10.0** have observations with a value of 0
 - Does this mean there was no particulate matter in the air?
 - More likely represents missing data, or for some reason PM not recorded
 - This also affects `Total_PM` and `Safe` variables
 - Table: example of this with `PM1.0`

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
0	697785	370.92	411095
0	697784	371.10	197567
0	594459	390.04	658137
0	594458	396.35	492275
0	594455	537.16	658135

+ Exploring the Data: Looking for Problems

- `proc univariate`
- **Temperature (degrees F)** has incorrect observations
 - -225, 131, 217, 255 are examples of incorrect temperatures!
 - This data will need to be accounted for
- **Humidity (%)** has incorrect observations
 - 1 observation of 0 % humidity
 - Several observations of 255 % humidity
 - Table: Humidity incorrect observations
- The sensors enter “+/-255” or “0” when there is an inability to read, sense, and save air quality data... these are values to look out for!

Extreme Observations			
Lowest		Highest	
Value	Obs	Value	Obs
0	461136	255	711704
11	394931	255	715364
11	394930	255	717961
11	394929	255	717982
11	394928	255	722066

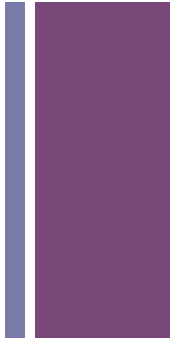
+ Exploring the Data



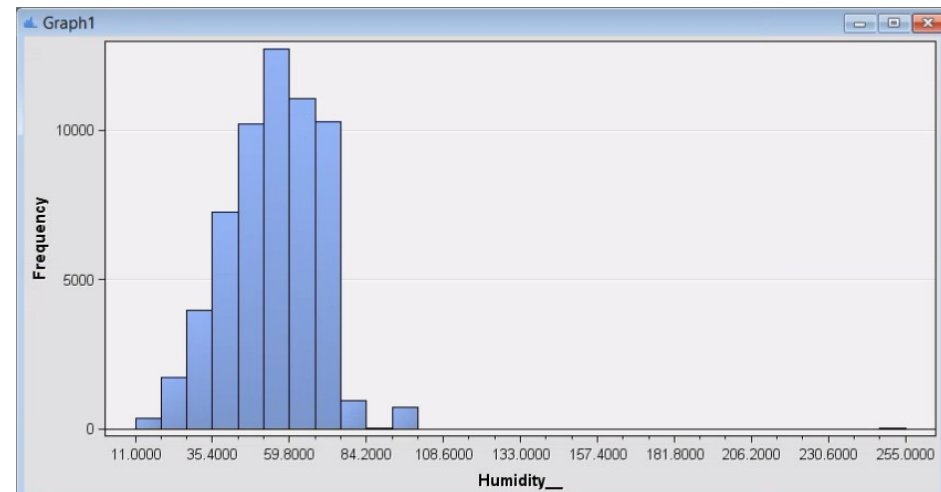
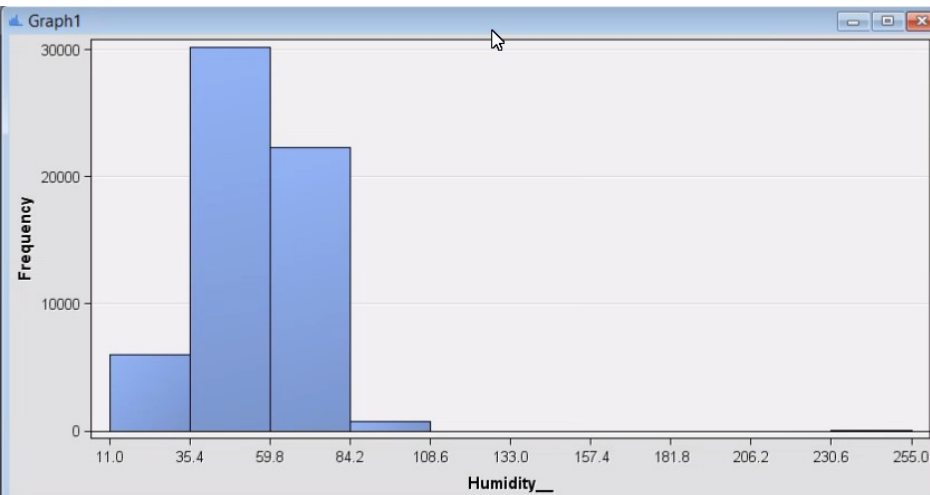
- Using SAS Enterprise Miner

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit
ADC	Input	Interval	No		No	.	.
Humidity	Input	Interval	No		No	.	.
PM10_0_CF_ATT	Input	Interval	No		No	.	.
PM1_0_CF_ATM	Input	Interval	No		No	.	.
PM2_5_CF_ATM	Input	Interval	No		No	.	.
Safe	Target	Binary	No		No	.	.
Temperature_F	Input	Interval	No		No	.	.
Total_PM	Input	Interval	No		No	.	.
UptimeMinutes	Input	Interval	No		No	.	.
created_at	Time ID	Interval	No		No	.	.
created_at_Strin	Time ID	Nominal	No		No	.	.
entry_jd	ID	Nominal	No		No	.	.
location	Input	Nominal	No		No	.	.
location_jd	ID	Nominal	No		No	.	.

+ Exploring the Data: Histograms

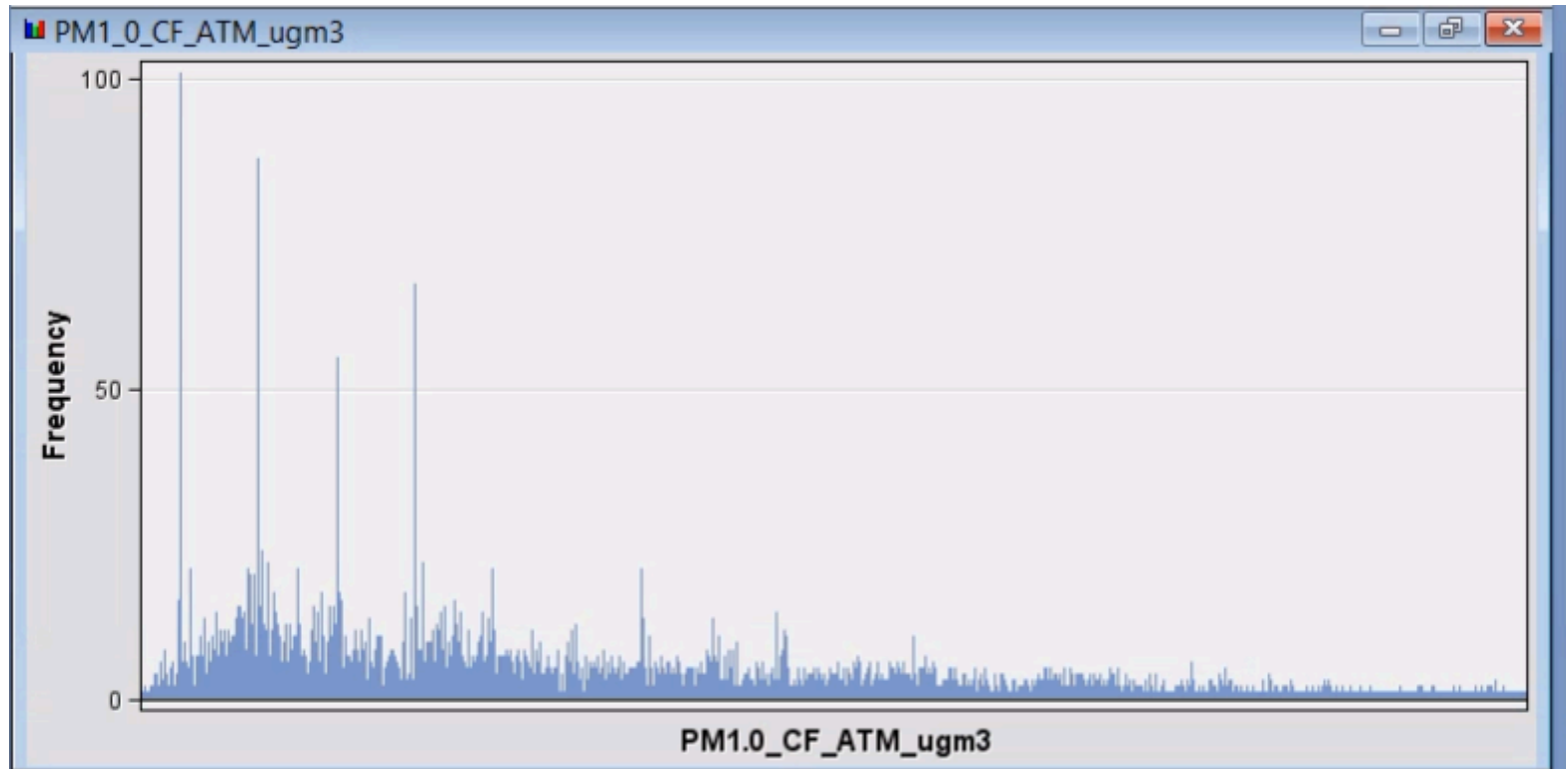


- I used the maximum of 30 bins for each independent variable
- This gave me the best, detailed view of each histogram
- Histogram: Humidity with 10 vs. 30 bins – better picture w/ more bins

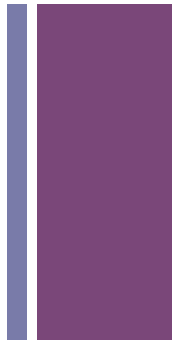


+ Exploring the Data: Histograms

- We have some variables that are skewed and will need to be transformed and normalized before we use them to create a model
 - PM 1.0, 2.5, and 10.0 variables are all right-skewed



+ Exploring the Data: Statistics



- Is the mean close to the median?
 - Most within 10 or less
- Skewness
 - Temperature is left skewed
 - ADC, PM variables, Safe are right skewed
- Kurtosis
 - ADC, PM variables have the most positive kurtosis, highly peaked

The SAS System

The MEANS Procedure

Variable	Label	Mean	Median	Mode	Minimum	Maximum	Range	Lower Quartile	Upper Quartile	Std Dev	Skewness	Kurtosis
ADC	ADC	-78.8488053	-79.0000000	-79.0000000	-96.0000000	31.0000000	127.0000000	-82.0000000	-77.0000000	11.0538267	8.6416324	83.5066683
created_at	created_at	1853521078	1853231679	1838302661	1838160035	1869696973	31535938.00	1845726272	1860963829	9027404.41	0.0798503	-1.1633741
Humidity	Humidity	54.7817903	55.0000000	55.0000000	0	255.0000000	255.0000000	45.0000000	66.0000000	14.4731152	0.0733987	1.6788523
PM10_0_CF_ATM_ugm3	PM10.0_CF_ATM_ugm3	13.4324544	9.6100000	1.0000000	0	2820.15	2820.15	4.1500000	18.0900000	14.7648911	20.3137780	2482.05
PM1_0_CF_ATM_ugm3	PM1.0_CF_ATM_ugm3	8.7846969	6.6500000	1.0000000	0	537.1600000	537.1600000	2.7200000	12.3100000	8.5494006	4.6613262	107.2969942
PM2_5_CF_ATM_ugm3	PM2.5_CF_ATM_ugm3	11.6923781	8.5800000	1.0000000	0	1099.29	1099.29	3.6000000	16.0700000	11.9973621	7.7419792	327.6363976
Total_PM	Total_PM	33.8895016	24.9000000	3.0000000	0	4135.77	4135.77	10.4700000	46.4800000	35.0063925	9.3574939	534.2396628
Temperature_F	Temperature_F	53.2400189	58.0000000	75.0000000	-225.0000000	217.0000000	442.0000000	36.0000000	76.0000000	39.1681749	-4.3103256	28.3159260
Safe	Safe	0.1033848	0	0	0	1.0000000	1.0000000	0	0	0.3044610	2.6053668	4.7879494
location_id	location_id	0.5047383	1.0000000	1.0000000	0	1.0000000	1.0000000	0	1.0000000	0.4999779	-0.0189540	-1.9996463

+ Cleaning the Data: Replacing Missing Data

- No data is missing from original data set, do not need to impute

Sample Statistics

Obs #	Variable Name	Label	Type	Percent ...	Minimum	Maximum	Mean	Number o...	Mode Per...	Mode
1	location	location	
2	created_at_String	created_at_...	CLASS	0	.	.	.	128+	0.775194	2018-04-01...
3	ADC	ADC	VAR	0	-83	-69	-76.4352.	.	.	.
4	Humidity	Humidity	VAR	0	27	73	47.7165.	.	.	.
5	PM10_0_CF_ATM_ugm3	PM10.0_CF...	VAR	0	0.7	54.72	11.00281.	.	.	.
6	PM1_0_CF_ATM_ugm3	PM1.0_CF...	VAR	0	0	35	7.054203.	.	.	.
7	PM2_5_CF_ATM_ugm3	PM2.5_CF...	VAR	0	0.25	47.25	9.453543.	.	.	.
8	Safe	Safe	VAR	0	0	1	0.1335.	.	.	.
9	Temperature_F	Temperatur...	VAR	0	21	50	35.77367.	.	.	.
10	Total_PM	Total_PM	VAR	0	1.77	133.82	27.51056.	.	.	.
11	UptimeMinutes	UptimeMin...	VAR	0	0	1799	525.9792.	.	.	.
12	created_at	created_at	VAR	0	1.8382E9	1.8386E9	1.8384E9.	.	.	.
13	entry_id	entry_id	VAR	0	527938	533937	530937.5.	.	.	.
14	location_id	location_id	VAR	0	1	1	1.	.	.	.

+ Cleaning the Data: Replacing Bad Data



- Using SAS Enterprise Miner – **Replacement Node**
 - Default Limits Method: None
 - Replacement Value: Missing – puts a missing indicator instead of a value
- **User-Specified Training** based on errors I uncovered
 - Humidity cannot be over 100% or under 0%
 - PM readings of 0 are an error
 - Created a reasonable Temperature range for Iowa

+ Cleaning the Data: Replacing Bad Data

Limits and Replacement Values for Interval Variables

Variable	Replace Variable	Lower limit	Lower Replacement Value	Upper Limit	Upper Replacement Value
Humidity	REP_Humidity	0.01	.	100	.
Temperature_F	REP_Temperature_F	-50.00	.	120	.
Total_PM	REP_Total_PM	0.01	.	.	.

```
*-----*  
* Report Output  
*-----*
```

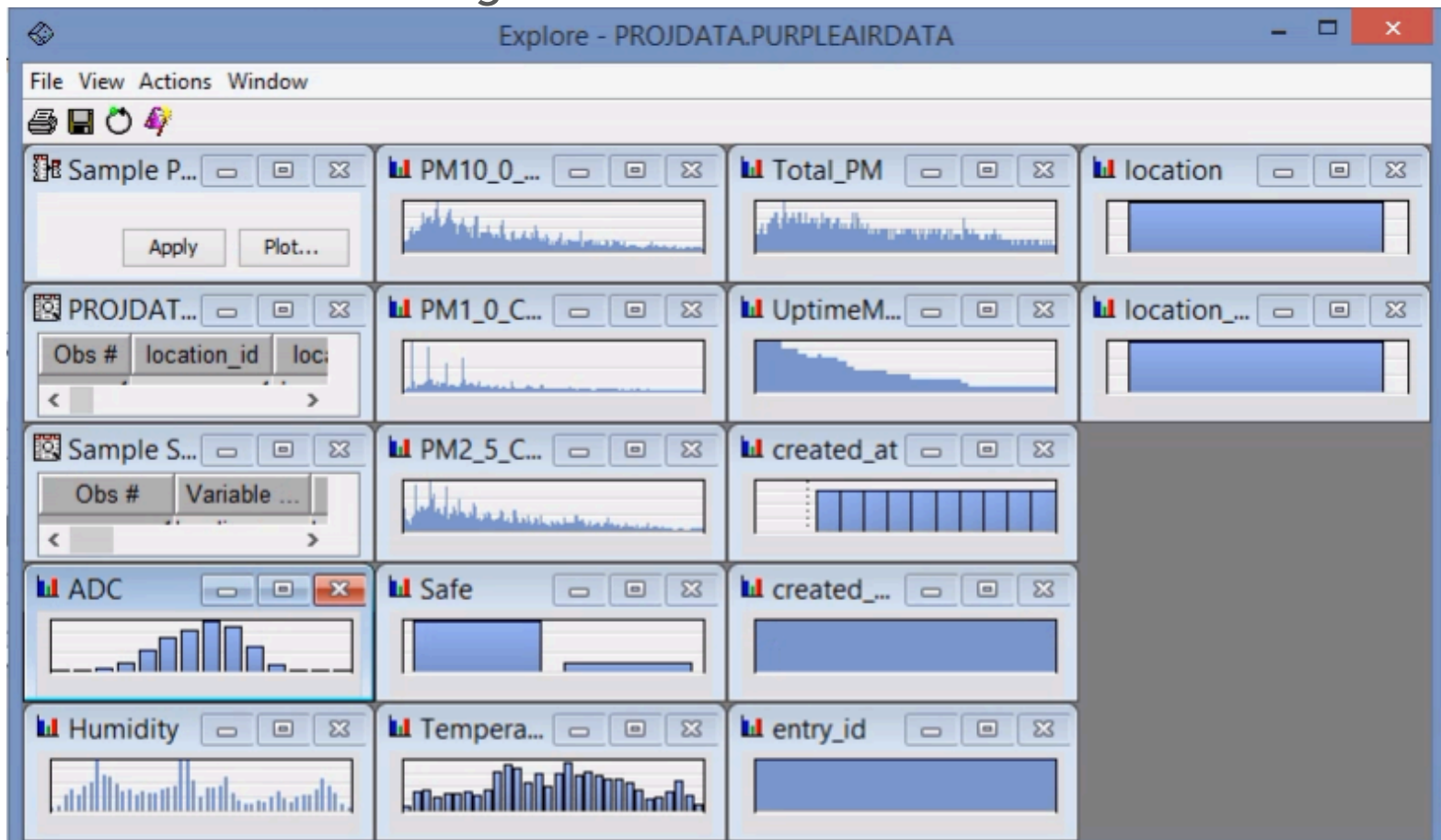
Replacement Counts

Obs	Variable	Label	Role	Train
1	Humidity	Humidity	INPUT	31
2	Temperature_F	Temperature_F	INPUT	8735
3	Total_PM	Total_PM	INPUT	52



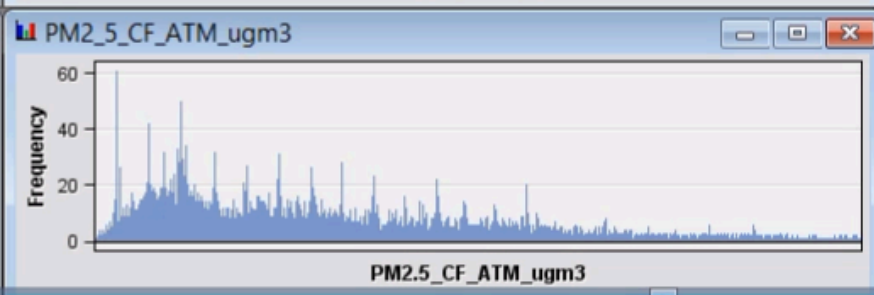
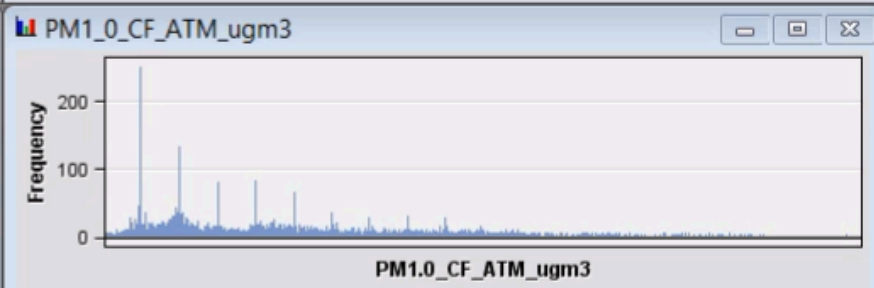
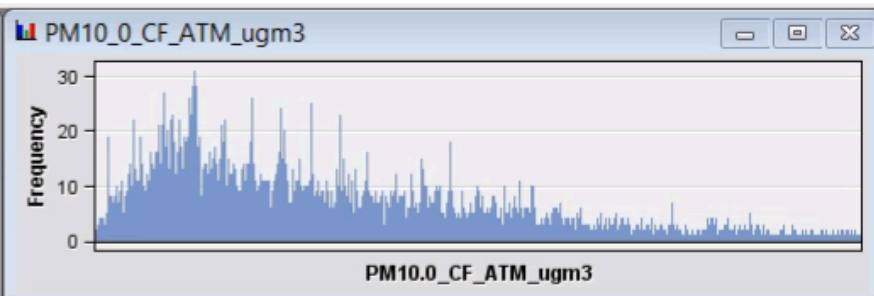
Transforming the Data: Fixing Skew

- We saw from the Histograms that several variables are non-normal and right-skewed



+ Transforming the Data: Fixing Skew

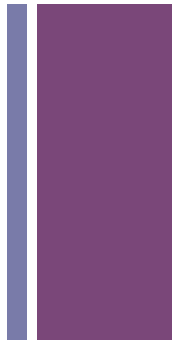
- Using SAS Enterprise Miner – **Transform Variables Node**
 - Any variables appearing right skewed, changed Method to Log



Name	Method	Number of Bins	Role	Level
ADC	Default	4	Input	Interval
Humidity	Default	4	Rejected	Interval
PM10_0_CF_ATM	Log	4	Input	Interval
PM1_0_CF_ATM	Log	4	Input	Interval
PM2_5_CF_ATM	Log	4	Input	Interval
REP_Humidity	Default	4	Input	Interval
REP_Temperature	Default	4	Input	Interval
REP_Total_PM	Default	4	Input	Interval
Safe	Default	4	Target	Binary
Temperature_F	Default	4	Rejected	Interval
Total_PM	Default	4	Rejected	Interval
UptimeMinutes	Default	4	Input	Interval
created_at_String	Default	4	Input	Nominal
location	Default	4	Input	Nominal



Transforming the Data: Fixing Skew

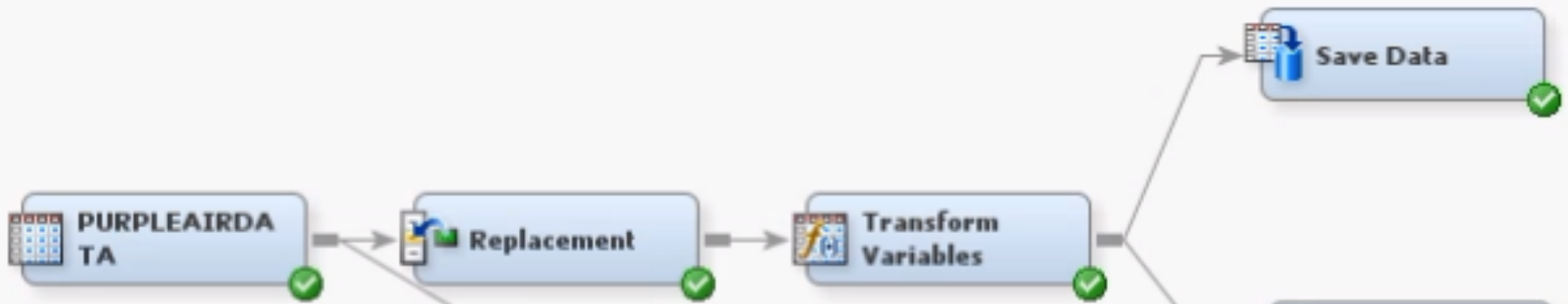


- Using SAS Enterprise Miner – **Transform Variables Node**
 - Any variables appearing right skewed, changed Method to Log
 - Normalized these variables
 - Mean and Standard Deviation
 - Skewness: now slightly left skewed, close to 0
 - Kurtosis: now slightly negative, close to 0, less peaked

Source	Method	Variable Name	Formula	Number of Levels	Non Missing	Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Label
Input	Original	PM10_0_C...		.	725687	0	0	2820.15	13.43245	14.76489	20.31378	2482.048	PM10.0_CF...
Input	Original	PM1_0_CF...		.	725687	0	0	537.16	8.764696	8.549401	4.661326	107.296	PM1.0_CF_...
Input	Original	PM2_5_CF...		.	725687	0	0	1099.29	11.69238	11.99736	7.741979	327.6364	PM2.5_CF_...
Output	Computed	LOG_PM10...	log(PM10_...	.	725687	0	0	7.9449	2.277841	0.929278	-0.25249	-0.43874	Transforme...
Output	Computed	LOG_PM1...	log(PM1_0...	.	725687	0	0	6.288156	1.933897	0.87491	-0.27018	-0.59685	Transforme...
Output	Computed	LOG_PM2...	log(PM2_5...	.	725687	0	0	7.003329	2.163673	0.916491	-0.27711	-0.49953	Transforme...

+ SAS Enterprise Miner Diagram

- We will use this prepared data for future models and analysis
- `em_save_train.sas7bdat` prepared SAS data file is attached to this submission



+ Next... Analyzing the Data

- Familiarize ourselves with the data
- Explore correlations & relationships between variables
- Create Models
 - Regression in Base SAS
 - Regression in SAS EM
 - Decision Tree (all variables included)
 - Decision Tree 2 (selected variables)
- Compare models
- Analyze results in context



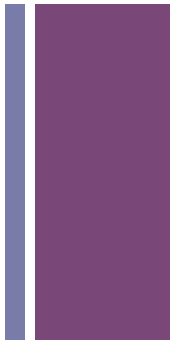
Data Exploration: Familiarization

- How does air quality vary with the time of day?
- Trends higher around rush hour ~ 8-9am and 5-6pm
- Spikes Common – Affected by immediate surroundings

1	downtown	2018-04-02 17:00:23 UTC	529782	9.57	12.32	13.61	92	-80	42	37	35.5	0
1	downtown	2018-04-02 17:01:41 UTC	529783	14.58	18.19	20.35	93	-76	42	38	53.12	1
1	downtown	2018-04-02 17:03:01 UTC	529784	15.16	18.88	20.81	95	-75	42	38	54.85	1
1	downtown	2018-04-02 17:04:21 UTC	529785	8.05	10.23	11.55	96	-75	42	37	29.83	0
1	downtown	2018-04-02 17:05:42 UTC	529786	6.69	8.36	9.76	97	-74	41	37	24.81	0
1	downtown	2018-04-02 17:07:01 UTC	529787	7.02	8.75	9.2	99	-79	42	37	24.97	0
1	downtown	2018-04-02 17:08:22 UTC	529788	6.71	8.49	9.15	100	-77	42	37	24.35	0
1	downtown	2018-04-02 17:09:41 UTC	529789	7.89	10.2	11.55	101	-74	42	37	29.64	0
1	downtown	2018-04-02 17:11:01 UTC	529790	8.23	10.56	11.33	103	-77	42	37	30.12	0
1	downtown	2018-04-02 17:12:21 UTC	529791	7.68	10.11	10.77	104	-74	42	37	28.56	0
1	downtown	2018-04-02 17:13:41 UTC	529792	7.09	9.51	11.95	105	-77	43	37	28.55	0
1	downtown	2018-04-02 17:15:01 UTC	529793	7.67	9.93	11.74	107	-75	41	36	29.34	0
1	downtown	2018-04-02 17:16:21 UTC	529794	6.92	8.79	9.69	108	-74	42	37	25.4	0
1	downtown	2018-04-02 17:17:41 UTC	529795	7.07	8.77	10.27	109	-75	42	37	26.11	0
1	downtown	2018-04-02 17:19:01 UTC	529796	6.95	8.43	9.19	111	-78	42	37	24.57	0
1	downtown	2018-04-02 17:20:21 UTC	529797	7.02	8.41	10.91	112	-77	42	37	26.34	0
1	downtown	2018-04-02 17:21:41 UTC	529798	5.07	7.29	7.67	113	-73	42	37	20.03	0
1	downtown	2018-04-02 17:23:02 UTC	529799	5.5	6.41	7.11	115	-74	42	36	19.02	0
1	downtown	2018-04-02 17:24:21 UTC	529800	6	7.29	8.63	116	-75	42	36	21.92	0
1	downtown	2018-04-02 17:25:42 UTC	529801	9.16	11.16	12.12	117	-75	42	36	32.44	0
1	downtown	2018-04-02 17:27:01 UTC	529802	7	9.21	10.74	119	-78	42	36	26.95	0
1	downtown	2018-04-02 17:28:22 UTC	529803	7.41	10.24	10.95	120	-76	42	36	28.6	0
1	downtown	2018-04-02 17:29:42 UTC	529804	8.3	10.14	11.91	121	-77	42	36	30.35	0
1	downtown	2018-04-02 17:31:02 UTC	529805	7.07	9.2	10.05	123	-76	42	36	26.32	0
1	downtown	2018-04-02 17:32:21 UTC	529806	7.43	9.31	11.52	124	-77	43	36	28.26	0
1	downtown	2018-04-02 17:33:41 UTC	529807	7.43	9.14	9.52	125	-79	43	36	26.09	0
1	downtown	2018-04-02 17:35:01 UTC	529808	22.6	28.33	32.44	127	-74	42	36	83.37	1
1	downtown	2018-04-02 17:36:21 UTC	529809	6.51	7.82	8.92	128	-78	42	36	23.25	0
1	downtown	2018-04-02 17:37:41 UTC	529810	4.52	6.5	8.95	130	-72	42	36	18.97	0



Data Exploration: Familiarization



■ Worst air quality?

4,135.77 Total_PM
downtown
5/5/18
3:35pm

■ Best air quality?

0.02 Total_PM
rural
10/11/18
from 8-11am

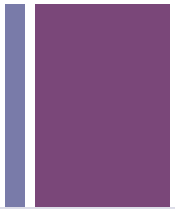
location_id	N Obs	Variable	Label	Mean	Median	Mode	Minimum	Maximum	Range
0	359405	ADC	ADC	-81.6355922	-81.0000000	-81.0000000	-96.0000000	31.0000000	127.0000000
		created_at	created_at	1853322988	1853225050	1854403200	1838160035	1869695922	31535887.00
		Humidity	Humidity	53.8394754	54.0000000	56.0000000	0	255.0000000	255.0000000
		PM10_0_CF_ATM_ugm3	PM10.0_CF_ATM_ugm3	11.6156655	8.3300000	1.0000000	0	1239.30	1239.30
		PM1_0_CF_ATM_ugm3	PM1.0_CF_ATM_ugm3	7.6494752	5.8200000	1.0000000	0	537.1600000	537.1600000
		PM2_5_CF_ATM_ugm3	PM2.5_CF_ATM_ugm3	10.1670129	7.4500000	1.0000000	0	917.6200000	917.6200000
		Total_PM	Total_PM	29.4321536	21.6400000	3.0000000	0	2694.08	2694.08
		Temperature_F	Temperature_F	50.8673876	61.0000000	-225.0000000	-225.0000000	217.0000000	442.0000000
		Safe	Safe	0	0	0	0	0	0
1	366282	ADC	ADC	-76.1143409	-78.0000000	-78.0000000	-96.0000000	31.0000000	127.0000000
		created_at	created_at	1853715450	1853238359	1845072000	1838160056	1869695973	31535917.00
		Humidity	Humidity	55.7043435	56.0000000	55.0000000	14.0000000	255.0000000	241.0000000
		PM10_0_CF_ATM_ugm3	PM10.0_CF_ATM_ugm3	15.2151329	11.0000000	1.0000000	0	2820.15	2820.15
		PM1_0_CF_ATM_ugm3	PM1.0_CF_ATM_ugm3	9.8589781	7.6000000	1.0000000	0	371.1000000	371.1000000
		PM2_5_CF_ATM_ugm3	PM2.5_CF_ATM_ugm3	13.1891042	9.8100000	1.0000000	0	1099.29	1099.29
		Total_PM	Total_PM	38.2631622	28.4700000	3.0000000	0	4135.77	4135.77
		Temperature_F	Temperature_F	55.5628925	56.0000000	74.0000000	-220.0000000	101.0000000	321.0000000
		Safe	Safe	0.2048285	0	0	0	1.0000000	1.0000000

+ Data Exploration: Correlation

- **How are PM 1.0, 2.5, 10.0, and Total related?**
- They are correlated with one another with a Prob $< |r|$ of $<.0001$
- If there is more of one PM size in the air, there is likely more of the other PM sizes in the air as well
 - Makes sense with what we know about PM behavior
- Total_PM and Safe are a function of PM 1.0, 2.5, 10.0
- These variables are redundant
- Will need to exclude them from our model

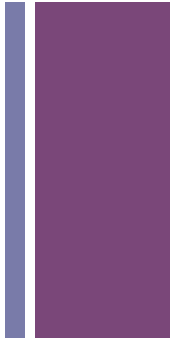


Data Exploration: Correlation



Pearson Correlation Coefficients Prob > r under H0: Rho=0 Number of Observations										
	ADC	created_at	Humidity	PM10_0_CF_ATM_ugm3	PM1_0_CF_ATM_ugm3	PM2_5_CF_ATM_ugm3	Total_PM	Temperature_F	Safe	
ADC	1.00000	0.20496	-0.01671	0.04433	0.04101	0.04266	0.04333	-0.10480	0.06871	
ADC		<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	725687	716531	725687
created_at	0.20496	1.00000	-0.05904	0.08579	0.05289	0.06965	0.07297	-0.29873	-0.02731	
created_at			<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	725687	716531	725687
Humidity	-0.01671	-0.05904	1.00000	0.16129	0.16256	0.16804	0.16534	-0.30080	0.12300	
Humidity		<.0001		<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001
	716531	716531	716531	716531	716531	716531	716531	716531	716531	716531
PM10_0_CF_ATM_ugm3	0.04433	0.08579	0.16129	1.00000	0.94786	0.98419	0.99000	0.02910	0.45590	
PM10_0_CF_ATM_ugm3		<.0001	<.0001		<.0001	<.0001	<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	716531	725687	725687
PM1_0_CF_ATM_ugm3	0.04101	0.05289	0.16256	0.94786	1.00000	0.98685	0.98165	0.05658	0.47692	
PM1_0_CF_ATM_ugm3		<.0001	<.0001	<.0001		<.0001	<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	716531	725687	725687
PM2_5_CF_ATM_ugm3	0.04266	0.06965	0.16804	0.98419	0.98685	1.00000	0.99830	0.04461	0.47516	
PM2_5_CF_ATM_ugm3		<.0001	<.0001	<.0001	<.0001		<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	716531	725687	725687
Total_PM	0.04333	0.07297	0.16534	0.99000	0.98165	0.99830	1.00000	0.04138	0.47161	
Total_PM		<.0001	<.0001	<.0001	<.0001	<.0001		<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	716531	725687	725687
Temperature_F	-0.10480	-0.29873	-0.30080	0.02910	0.05658	0.04461	0.04138	1.00000	0.04685	
Temperature_F		<.0001	<.0001	<.0001	<.0001	<.0001	<.0001		<.0001	<.0001
	716531	716531	716531	716531	716531	716531	716531	716531	716531	716531
Safe	0.06871	-0.02731	0.12300	0.45590	0.47692	0.47516	0.47161	0.04685	1.00000	
Safe		<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001		<.0001
	725687	725687	716531	725687	725687	725687	725687	716531	725687	725687

+ Data Exploration: Correlation



The SAS System

The CORR Procedure

1 With Variables:	Safe
8 Variables:	ADC created_at Humidity PM10_0_CF_ATM_ugm3 PM1_0_CF_ATM_ugm3 PM2_5_CF_ATM_ugm3 Total_PM Temperature_F

Pearson Correlation Coefficients
 Prob > |r| under H0: Rho=0
 Number of Observations

	ADC	created_at	Humidity	PM10_0_CF_ATM_ugm3	PM1_0_CF_ATM_ugm3	PM2_5_CF_ATM_ugm3	Total_PM	Temperature_F
Safe	0.06871	-0.02731	0.12300	0.45590	0.47692	0.47516	0.47161	0.04685
Safe	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001	<.0001
	725687	725687	716531	725687	725687	725687	725687	716531

+ Question Review



- **What affects air quality?**
- **Can we predict if air quality is safe using only data readily available on any basic weather report?**
(Temperature, Humidity, Date, Location, etc.)

+ Review of Variables

- Date and time air quality reading was collected
- Entry Identifier
- PM1.0 CF ATM ug/m³ (particulate matter 1 μm or less in diameter)
- PM2.5 CF ATM ug/m³ (particulate matter 2.5 μm or less in diameter)
- PM 10.0 CF ATM ug/m³ (particulate matter 10 μm or less in diameter)
- Total PM – a user-created variable of total PM 1.0, 2.5, and 10.0 in the air
- Safe – user-created variable of if air quality is safe or not (1=safe, 0=unsafe)
- Uptime Minutes
- Temperature (degrees Fahrenheit)
- Humidity (as a percent)
- Analog to Digital Converter (ADC) reading
- Location – user-created character variable
- Location Identifier – user created binary variable with 1=downtown, 0=rural

μm = micrometer

CF ATM = cubic foot of atmosphere

PM = particulate matter

+ Variable Roles



- Decision Variable: **Safe**
 - This will be my target variable
- Independent Variables:
 - **Humidity__**
 - **Temperature_F**
 - **Created_At**
 - **Location_ID**
 - These will be predictor variables

+ Modeling Method: Logistic



- We have a Binary Decision variable – perfect for Logistic Regression
- **Base SAS**
 - Examine Data, Variables, and Relationships between variables
 - Proc corr
 - Proc univariate
 - Proc means
 - Proc sgplot
 - Preliminary modeling
 - Proc logistic
 - Proc freq

+ Modeling Method: Logistic

- $\text{Safe} = -21.6428 + 17.9505 * \text{Location_ID} + 0.03565 * \text{Humidity} + 0.00457 * \text{Temperature_F}$
- Temperature and Humidity are significant

Model Convergence Status	
Convergence criterion (GCONV=1E-8) satisfied.	

Analysis of Maximum Likelihood Estimates					
Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	-21.6428	25.4753	0.7217	0.3956
location_id	1	17.9505	25.4753	0.4965	0.4810
Humidity	1	0.0365	0.000369	9782.0850	<.0001
Temperature_F	1	0.00457	0.000193	558.9021	<.0001

Model Fit Statistics		
Criterion	Intercept Only	Intercept and Covariates
AIC	480509.79	357114.90
SC	480521.27	357160.83
-2 Log L	480507.79	357108.90

Testing Global Null Hypothesis: BETA=0			
Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	123400.888	3	<.0001
Score	92037.4675	3	<.0001
Wald	11469.2296	3	<.0001

Odds Ratio Estimates			
Effect	Point Estimate	95% Wald Confidence Limits	
location_id	>999.999	<0.001	>999.999
Humidity	1.037	1.036	1.038
Temperature_F	1.005	1.004	1.005

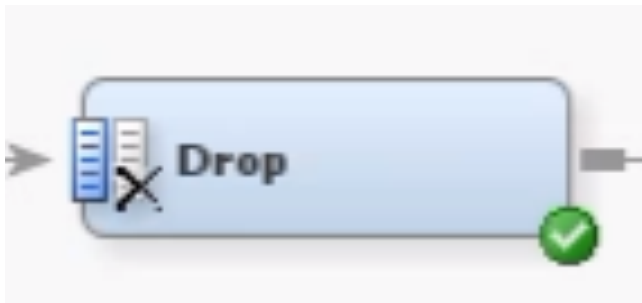
+ Modeling Method: Logistic



- It is 65 degrees F, 50% humidity, and you live 5 blocks from downtown.
- $\text{Safe} = -21.6428 + (17.9505 * 1) + (0.03565 * 50) + (0.00457 * 65)$
 $\text{Safe} = -1.61275$
- Will examine Logistic Regression Model closer in SAS EM
 - Prepared data
 - Training data

+ Data Exclusion

- **Drop Node** – put before every model
- Variables to be excluded
 - Redundant
 - Irrelevant

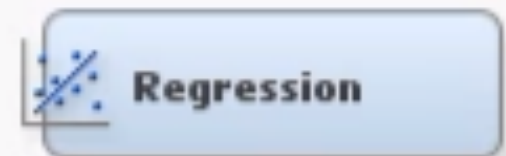


Name	Drop ▾	Role	Level
entry_id	Yes	ID	Nominal
ADC	Yes	Input	Interval
LOG_PM2_5_CF	Yes	Input	Interval
created_at_String	Yes	Input	Nominal
dataobs	Yes	ID	Interval
UptimeMinutes	Yes	Input	Interval
Total_PM	Yes	Rejected	Interval
LOG_PM1_0_CF	Yes	Input	Interval
REP_Total_PM	Yes	Input	Interval
LOG_PM10_0_CF	Yes	Input	Interval
REP_Temperature	Default	Input	Interval
location	Default	Input	Nominal
Safe	Default	Target	Binary
Temperature_F	Default	Rejected	Interval
Humidity	Default	Rejected	Interval
REP_Humidity	Default	Input	Interval
location_id	Default	ID	Binary
created_at	Default	Time ID	Interval

+ Modeling Method: Logistic



- We have a Binary Decision variable – perfect for Logistic Regression
- **SAS Enterprise Miner:**
 - Logistic Regression – **Regression Node**
 - Using Prepared Data Set - **Replace, Transform**
 - Partition - 67% train, 33% Validate
 - **Drop Node** – Remove irrelevant and unnecessary variables
 - **Safe** will be the Target Variable
 - Analyze Results
 - Analysis of Maximum Likelihood Estimates – Significance of each variable
 - Odds Ratios Estimates – Chances of Safe or U
 - AIC – Fit Statistic for the model



+ Modeling Method: Logistic

- **Analysis of Maximum Likelihood Estimates:**
Significance of each variable
- All variables used in the Selected model
- All variables PR < ChiSq is less than 0.05
 - All variables are significant to the model!
 - Estimate: Location has the most impact
 - Humidity and Temperature have small impact

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq	Standardized Estimate	Exp (Est)
Intercept	1	-10.6159	2.1204	25.07	<.0001		0.000
REP_Humidity	1	0.0364	0.000451	6507.59	<.0001	0.2729	1.037
REP_Temperature_F	1	0.00446	0.000236	356.59	<.0001	0.0597	1.004
location downtown	1	6.9347	2.1202	10.70	0.0011		999.000

+ Modeling Method: Logistic

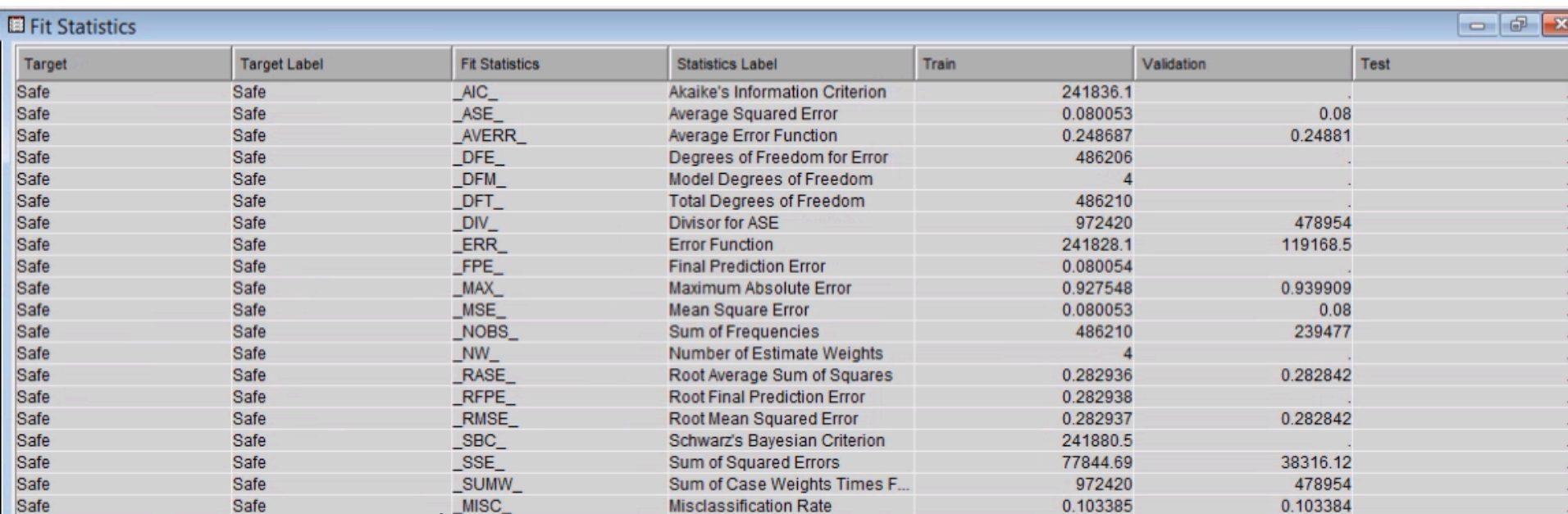
- **Odds Ratios Estimates:** Chances of Safe/Unsafe air quality
- For each additional % Humidity, the odds of Safe increase by 3.7%
- For each additional degree F, the odds of Safe increase by 0.4%
- Location is a determining factor: rural has higher chances of unsafe

Odds Ratio Estimates

Effect	Point Estimate
REP_Humidity	1.037
REP_Temperature_F	1.004
location	999.000

+ Modeling Method: Logistic

- **AIC:** Fit Statistic for the model
- 241,836.1 (better than Base SAS AIC of 357,114.50)

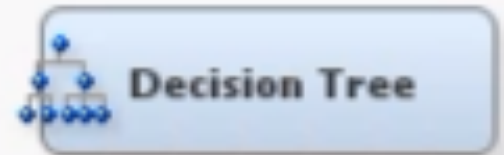


Target	Target Label	Fit Statistics	Statistics Label	Train	Validation	Test
Safe	Safe	_AIC_	Akaike's Information Criterion	241836.1	.	.
Safe	Safe	_ASE_	Average Squared Error	0.080053		0.08
Safe	Safe	_AVERR_	Average Error Function	0.248687		0.24881
Safe	Safe	_DFE_	Degrees of Freedom for Error	486206		.
Safe	Safe	_DFM_	Model Degrees of Freedom	4		.
Safe	Safe	_DFT_	Total Degrees of Freedom	486210		.
Safe	Safe	_DIV_	Divisor for ASE	972420		478954
Safe	Safe	_ERR_	Error Function	241828.1		119168.5
Safe	Safe	_FPE_	Final Prediction Error	0.080054		.
Safe	Safe	_MAX_	Maximum Absolute Error	0.927548		0.939909
Safe	Safe	_MSE_	Mean Square Error	0.080053		0.08
Safe	Safe	_NOBS_	Sum of Frequencies	486210		239477
Safe	Safe	_NW_	Number of Estimate Weights	4		.
Safe	Safe	_RASE_	Root Average Sum of Squares	0.282936		0.282842
Safe	Safe	_RFPE_	Root Final Prediction Error	0.282938		.
Safe	Safe	_RMSE_	Root Mean Squared Error	0.282937		0.282842
Safe	Safe	_SBC_	Schwarz's Bayesian Criterion	241880.5		.
Safe	Safe	_SSE_	Sum of Squared Errors	77844.69		38316.12
Safe	Safe	_SUMW_	Sum of Case Weights Times F...	972420		478954
Safe	Safe	_MISC_	Misclassification Rate	0.103385		0.103384

+ Modeling Method: Decision Tree



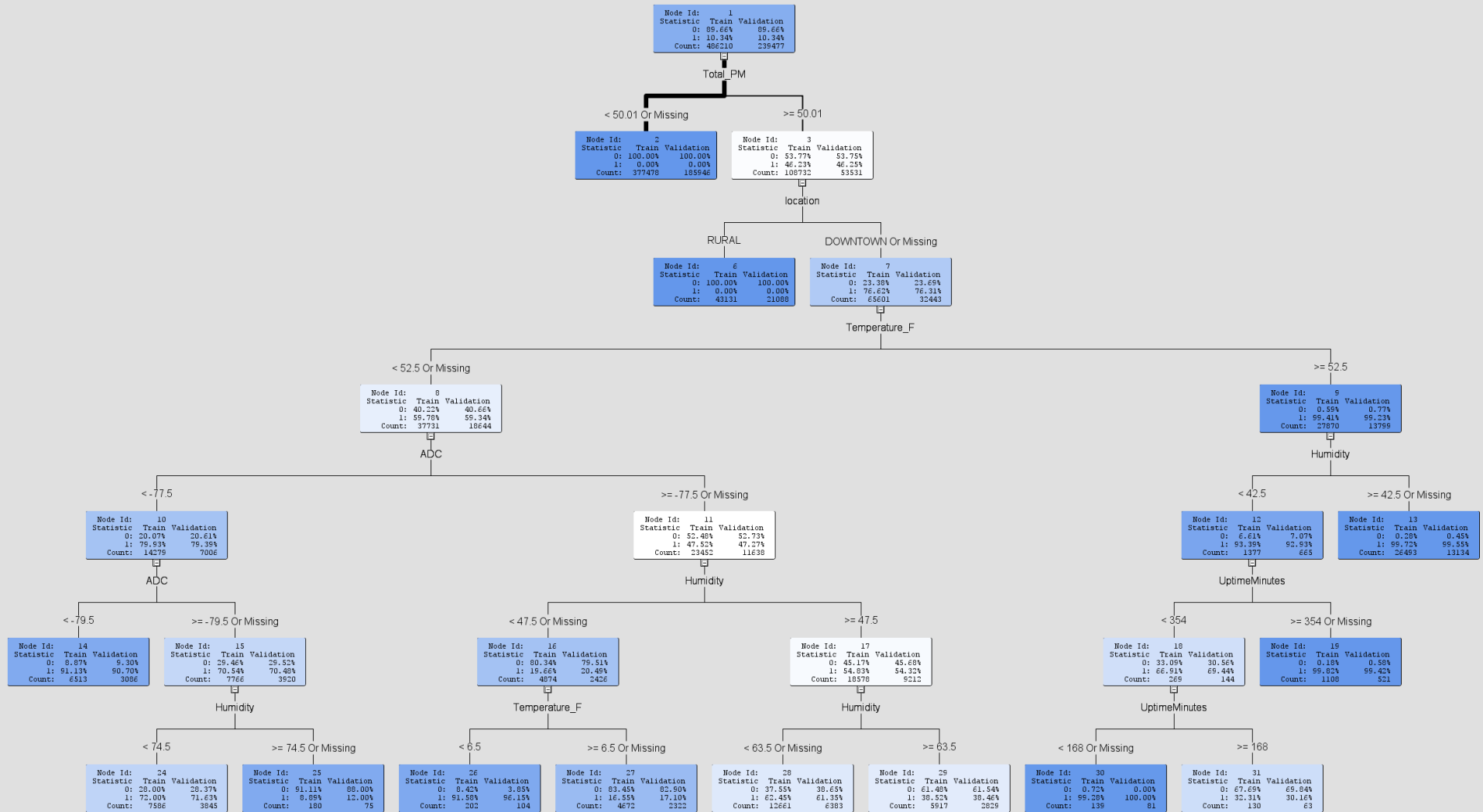
- A Decision Tree makes sense since we are trying to predict air quality and decide if it is safe or unsafe based on other factors/variables
- **SAS Enterprise Miner:**
 - Decision Tree - **Decision Tree Node**
 - Do NOT need to use Prepared Data Set
 - Will partition the data for training - **Data Partition Node**
 - 67% Train and 33% Validate
 - **Safe** will be the Target Variable
 - Analyze Results
 - Event Classification Table - True and False Negatives and Positives
 - Decision Tree - see what is the Root Node
 - AIC - Fit Statistic for the model





Modeling Method: Decision Tree

■ Curious... All Variables



+ Modeling Method: Decision Tree



Variable Importance

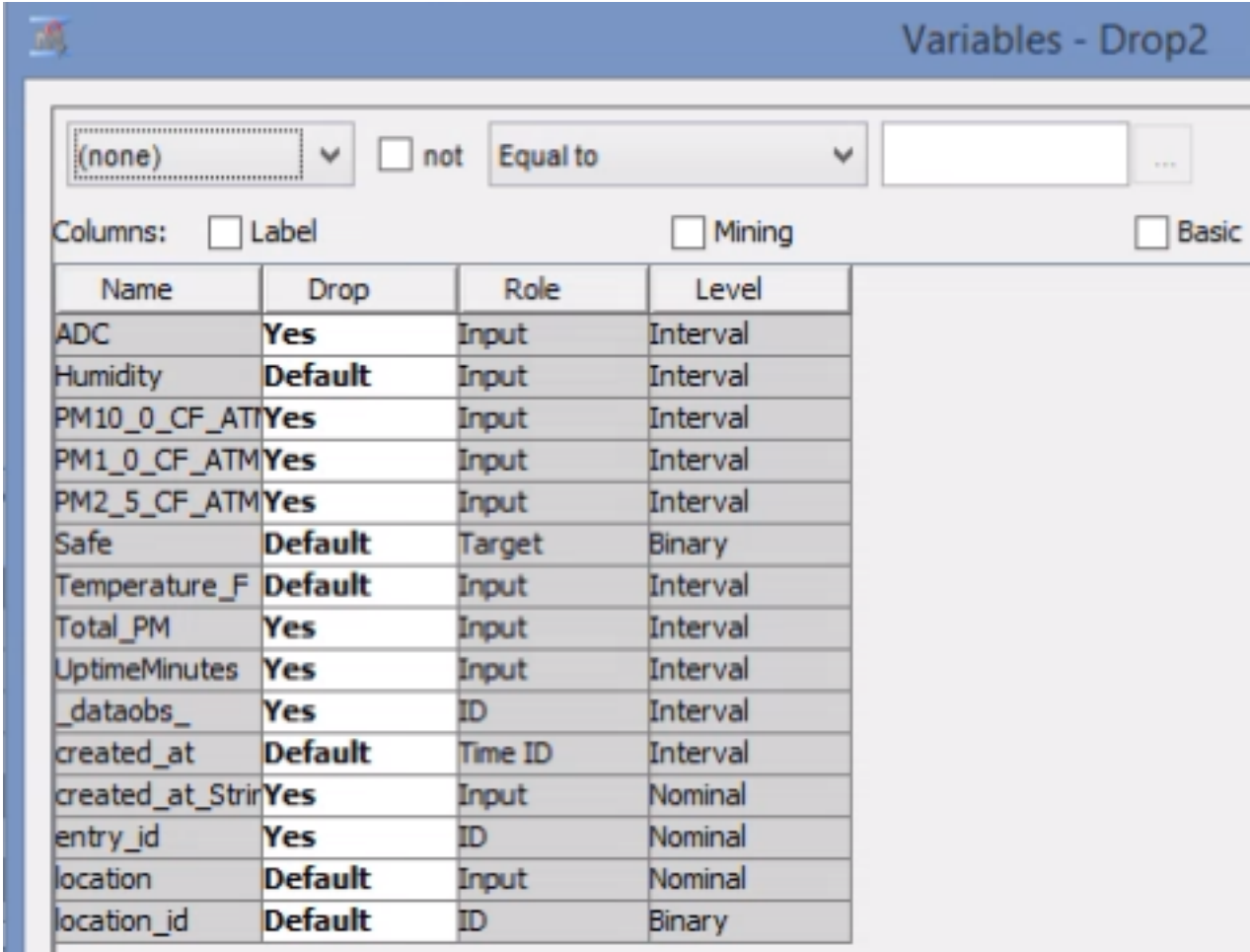
Variable Name	Label	Number of Splitting Rules	Importance	Validation Importance	Ratio of Validation to Training Importance
Total_PM	Total_PM	1	1.0000	1.0000	1.0000
location	location	1	0.9203	0.9149	0.9942
Temperature_F	Temperature_F	2	0.3815	0.3859	1.0115
ADC	ADC	2	0.2450	0.2423	0.9893
Humidity	Humidity	4	0.2085	0.1986	0.9524
UptimeMinutes	UptimeMinutes	2	0.0545	0.0555	1.0188

This is not helpful to answering our questions.

+ Modeling Method: Decision Tree

- This includes redundant and irrelevant
 - Variables not relevant to our questions!

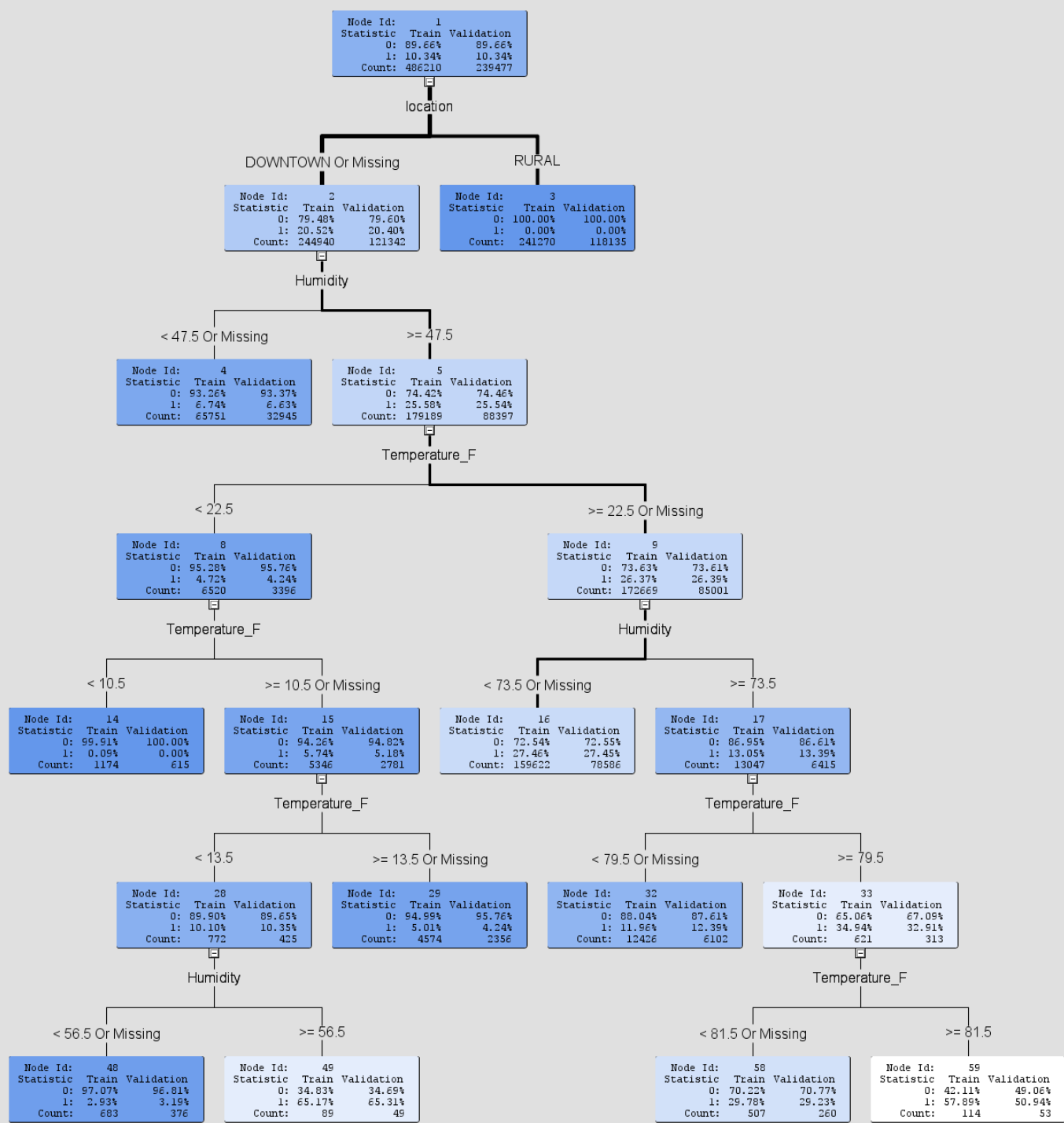
- **Drop Node** – removed variables



The screenshot shows a software interface titled "Variables - Drop2". At the top, there is a search bar with "(none)" selected, a checkbox for "not", and a dropdown menu set to "Equal to". Below this, there are checkboxes for "Columns:", "Label", "Mining", and "Basic". The main part of the interface is a table with the following data:

Name	Drop	Role	Level
ADC	Yes	Input	Interval
Humidity	Default	Input	Interval
PM10_0_CF_ATM	Yes	Input	Interval
PM1_0_CF_ATM	Yes	Input	Interval
PM2_5_CF_ATM	Yes	Input	Interval
Safe	Default	Target	Binary
Temperature_F	Default	Input	Interval
Total_PM	Yes	Input	Interval
UptimeMinutes	Yes	Input	Interval
dataobs	Yes	ID	Interval
created_at	Default	Time ID	Interval
created_at_Strin	Yes	Input	Nominal
entry_id	Yes	ID	Nominal
location	Default	Input	Nominal
location_id	Default	ID	Binary

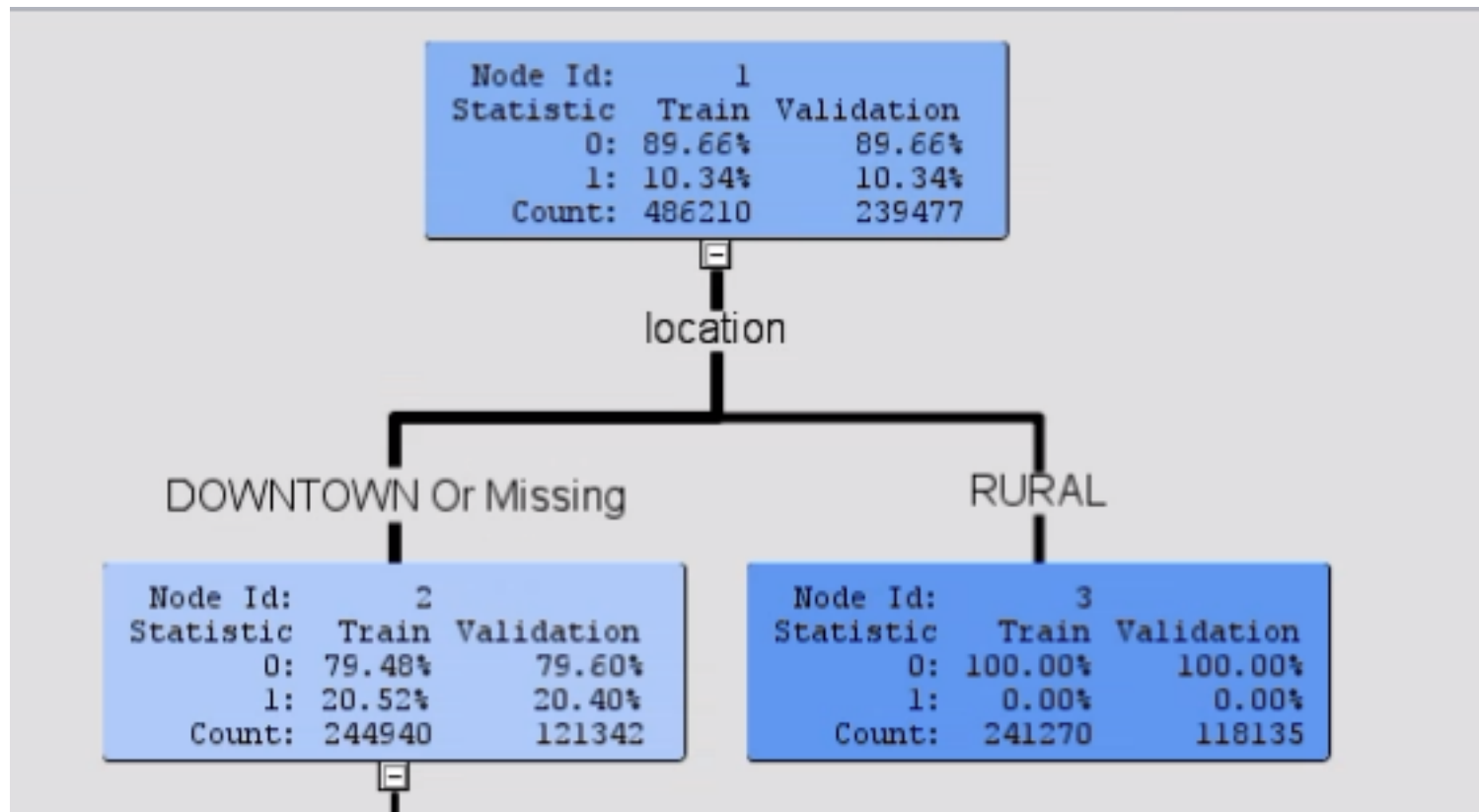
+ Modeling Method: Decision Tree 2



1 = Safe
0 = Unsafe

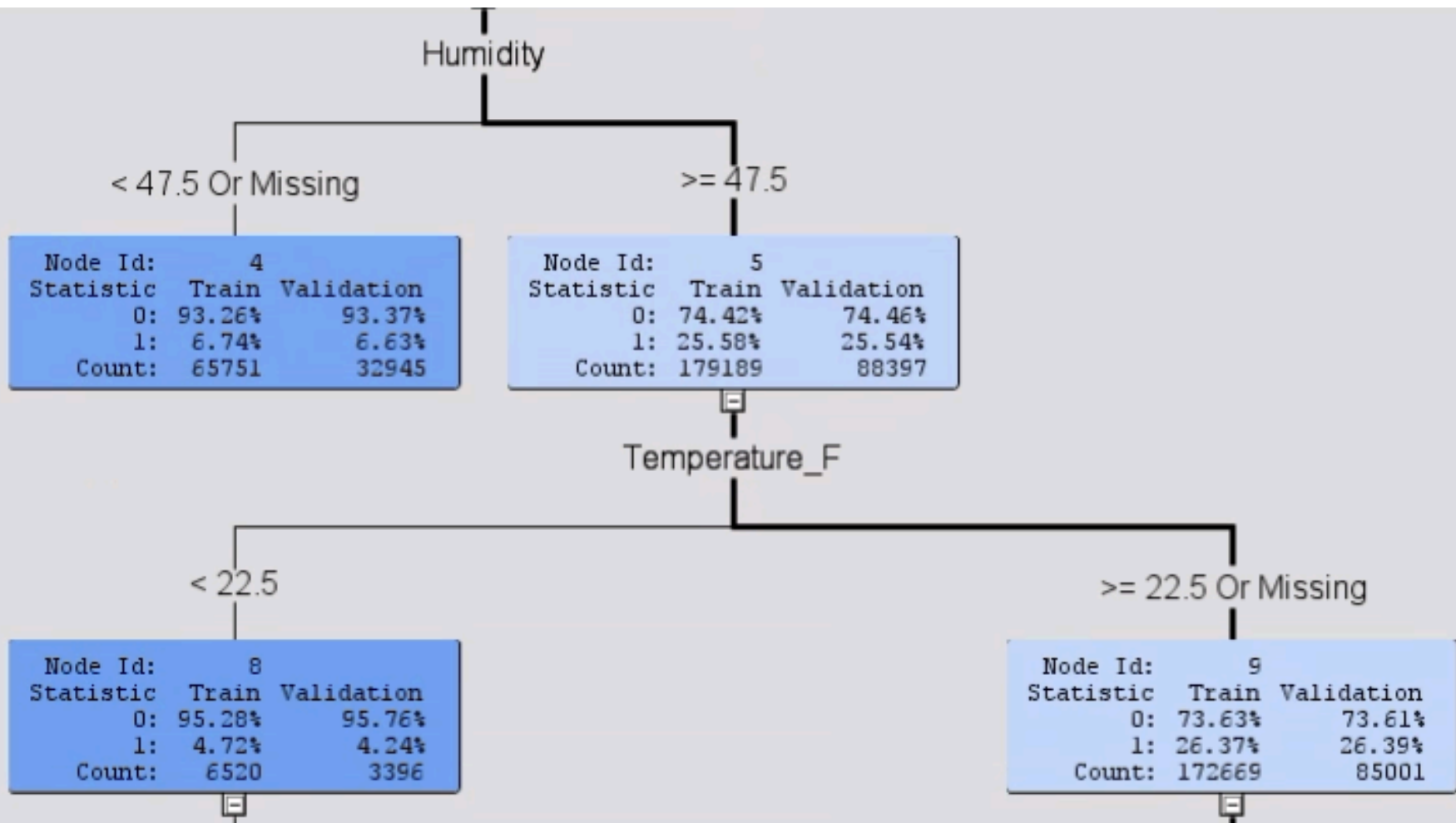
+ Modeling Method: Decision Tree 2

- **Root node** = Location
- RURAL has 0% Chance of Safe Air
DOWNTOWN has 20% Chance of Safe Air
 - Seems counterintuitive... unexpected results are still results! More later...



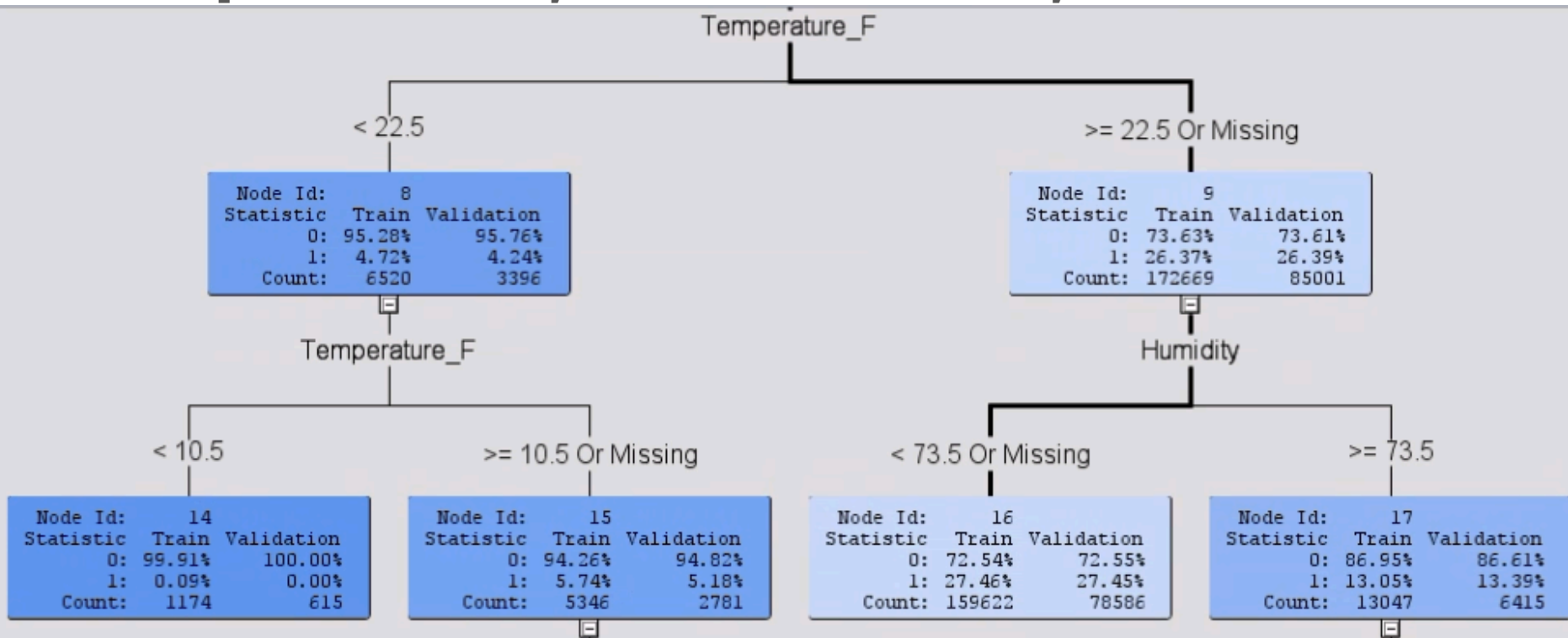
+ Modeling Method: Decision Tree 2

- Drier Air has more PM
- Humid Air has 25% & Dry Air has 6% chance of being Safe



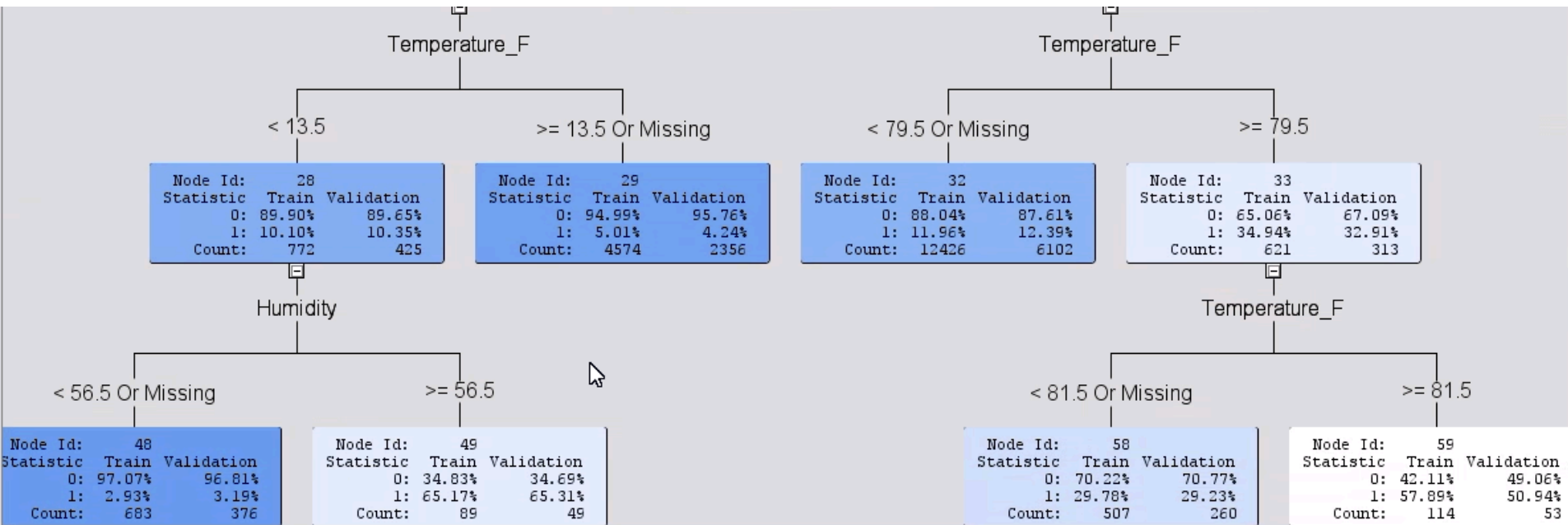
+ Modeling Method: Decision Tree 2

- Cold Air (10 degrees+ below freezing) 5% or less chance Safe
 - More Unsafe as it gets colder
- Warmer Air 25% Safe
 - depends on Humidity: 25-27% chance Safe if dry, 13% chance Safe if humid



+ Modeling Method: Decision Tree 2

- Further broken down by Temperature and Humidity
- In general:
 - Colder = More Likely to be Unsafe, Warm = More likely to be Safe
 - Dry = More Likely to be Unsafe, Humid = More Likely to be Safe
 - Cold and Dry is BAD! Warm and Humid is BETTER!



+ Modeling Method: Decision Tree 2

- Location is most important, followed by Humidity, then Temperature

Variable Importance

Variable Name	Label	Number of Splitting Rules	Importance	Validation Importance	Ratio of Validation to Training Importance
location	location	1	1.0000	1.0000	1.0000
Humidity	Humidity	3	0.6232	0.6310	1.0125
Temperature_F	Temperature_F	5	0.2568	0.2666	1.0378

+ Modeling Method: Decision Tree 2

■ Event Classification Table: True/False Negatives and Positives

Classification Table

Data Role=TRAIN Target Variable=Safe Target Label=Safe

Target	Outcome	Target Percentage	Outcome Percentage	Frequency Count	Total Percentage
0	0	89.6827	99.9819	435864	89.6452
1	0	10.3173	99.7533	50143	10.3130
0	1	38.9163	0.0181	79	0.0162
1	1	61.0837	0.2467	124	0.0255

Data Role=VALIDATE Target Variable=Safe Target Label=Safe

Target	Outcome	Target Percentage	Outcome Percentage	Frequency Count	Total Percentage
0	0	89.6819	99.9800	214676	89.6437
1	0	10.3181	99.7617	24699	10.3137
0	1	42.1569	0.0200	43	0.0180
1	1	57.8431	0.2383	59	0.0246

+ Modeling Method: Decision Tree 2

- Event Classification Table: True/False Negatives and Positives

Event Classification Table

Data Role=TRAIN Target=Safe Target Label=Safe

False Negative	True Negative	False Positive	True Positive
50143	435864	79	124

Data Role=VALIDATE Target=Safe Target Label=Safe

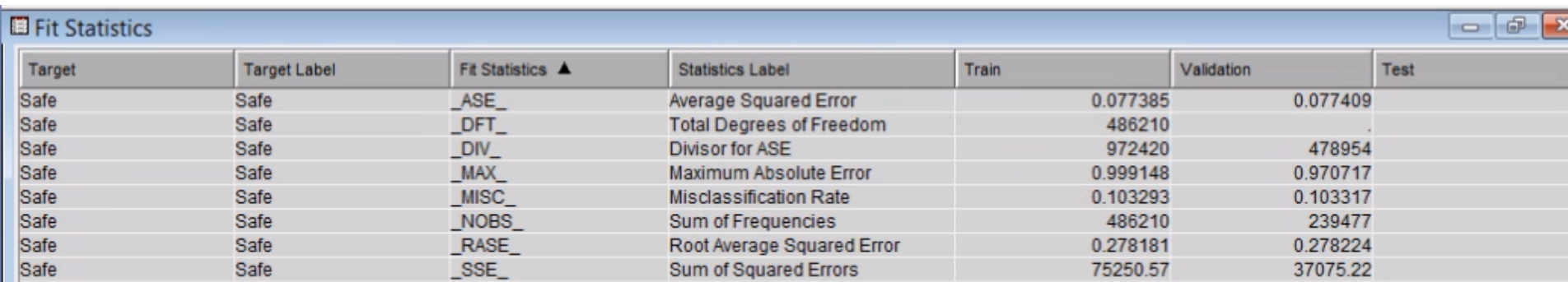
False Negative	True Negative	False Positive	True Positive
24699	214676	43	59

+ Modeling Method: Decision Tree 2

- Misclassification Rate for Validation Data Set: 10.3%
(False Positive + False Negatives) / Total
Only wrong 10.3% of the time
- Specificity for Validation Data Set: 99.98%
True Negative Predicted / Actual Negative
When it is false/Unsafe, how often does it predict false?
- Sensitivity for Validation Data Set: 0.2383%
True Positive Predicted / Actual Positive
When it is true/Safe, how often does it predict true?

+ Modeling Method: Decision Tree 2

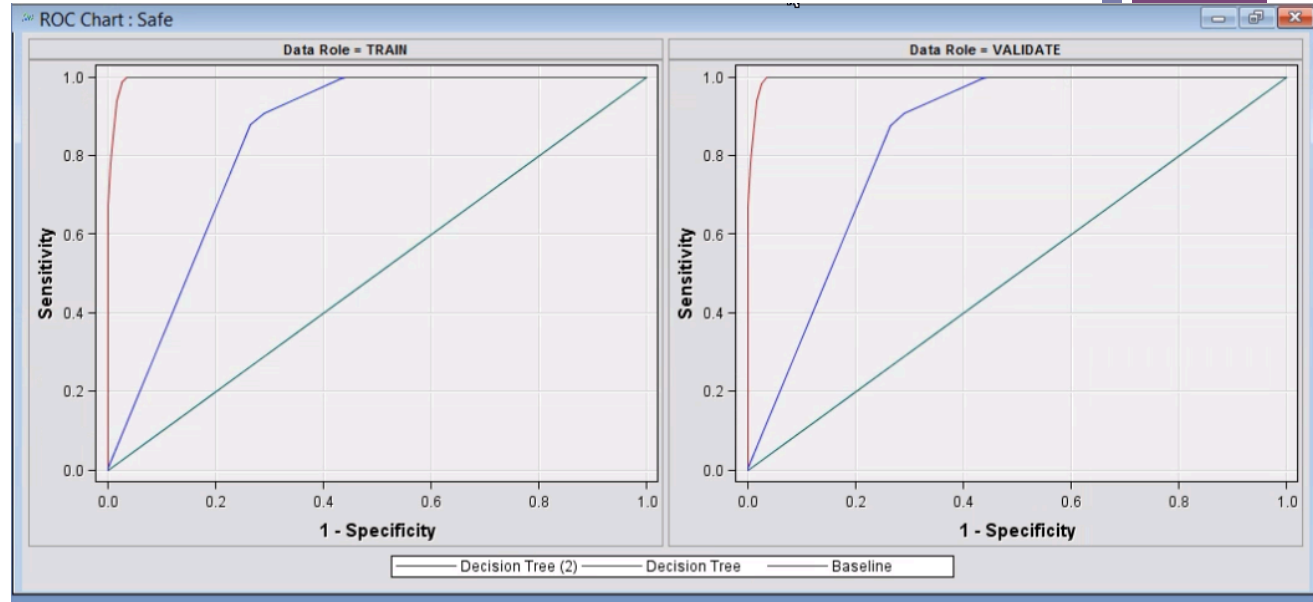
- Fit Statistics
- Average Square Error closer to 0 is better
- Similar in Train and Validation data sets:
Not over fitting or under fitting the model



Target	Target Label	Fit Statistics ▲	Statistics Label	Train	Validation	Test
Safe	Safe	_ASE_	Average Squared Error	0.077385	0.077409	
Safe	Safe	_DFT_	Total Degrees of Freedom	486210	.	.
Safe	Safe	_DIV_	Divisor for ASE	972420	478954	.
Safe	Safe	_MAX_	Maximum Absolute Error	0.999148	0.970717	.
Safe	Safe	_MISC_	Misclassification Rate	0.103293	0.103317	.
Safe	Safe	_NOBS_	Sum of Frequencies	486210	239477	.
Safe	Safe	_RASE_	Root Average Squared Error	0.278181	0.278224	.
Safe	Safe	_SSE_	Sum of Squared Errors	75250.57	37075.22	.

+ Compare Decision Trees

- Tree 1 (all variables) is better, but it doesn't make sense given our questions.
- Tree 2 is most helpful for someone who doesn't know air quality readings!



Fit Statistics

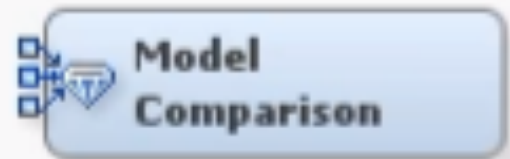
Model Selection based on Valid: Misclassification Rate (`_VMISC_`)

Selected Model	Model Node	Model Description	Valid: Misclassification Rate	Train: Average Squared Error	Train: Misclassification Rate	Valid: Average Squared Error
Y	Tree	Decision Tree	0.02266	0.014832	0.02193	0.015217
	Tree2	Decision Tree (2)	0.10332	0.077385	0.10329	0.077409

+ Modeling Method: Model Comparison

■ SAS Enterprise Miner:

- Compare Regression & Decision Tree 2 – **Model Comparison Node**
- Which does SAS EM think is better?
- Analyze Results
 - Fit Statistics – Which is Chosen, denoted with a Y
 - ROC Chart – Which has more under the curve?
 - Does this make sense?



+ Modeling Method: Model Comparison

- Fit Statistics – Decision Tree 2 is chosen
 - Slightly lower Misclassification Rates
 - Slightly lower Average Squared Errors

Fit Statistics

Model Selection based on Valid: Misclassification Rate (_VMISC_)

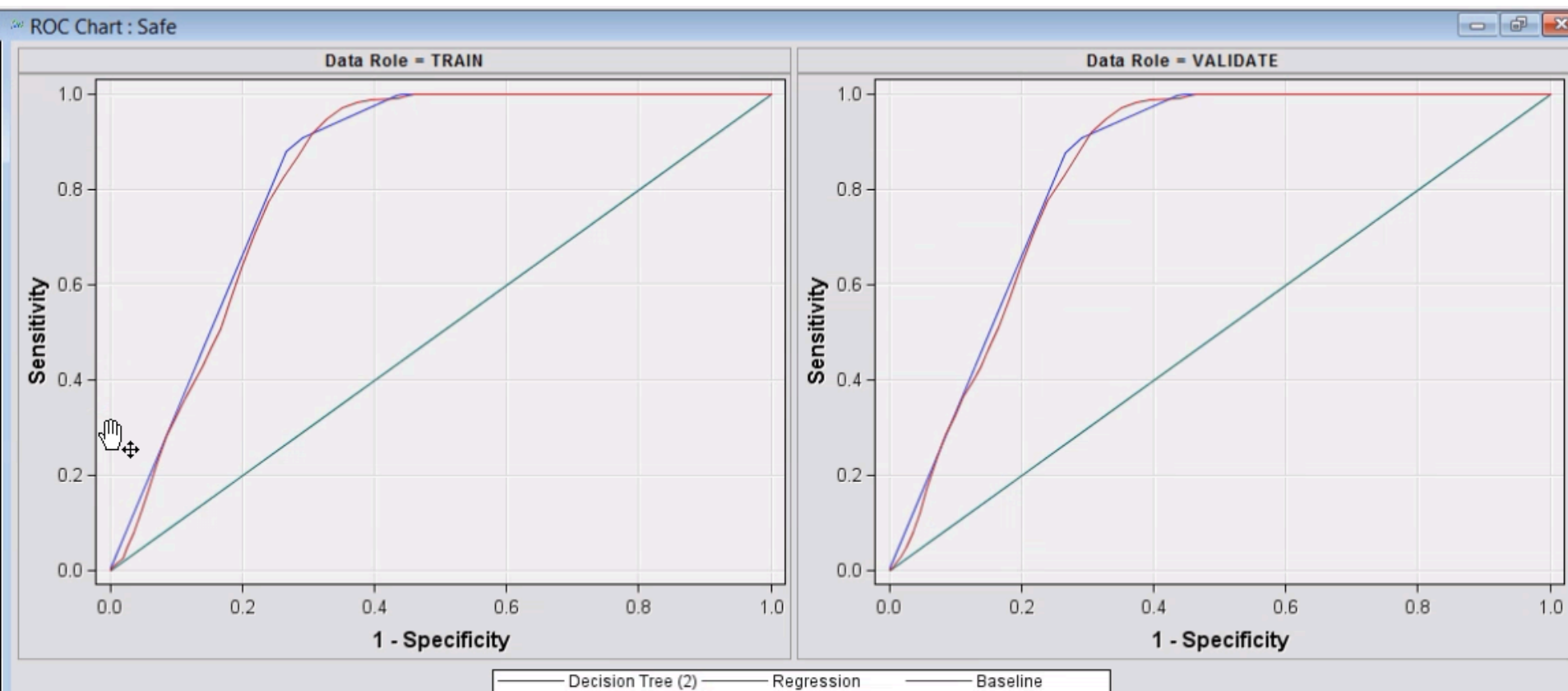
Selected Model	Model Node	Model Description	Valid: Misclassification Rate	Train: Average Squared Error	Train: Misclassification Rate	Valid: Average Squared Error
Y	Tree2	Decision Tree (2)	0.10332	0.077385	0.10329	0.077409
	Reg	Regression	0.10338	0.080053	0.10339	0.080000



Modeling Method: Model Comparison



- ROC Chart: Plots Sensitivity against Specificity
 - Decision Tree 2 closest to top left hand corner, more area under curve





Modeling Method: Model Comparison

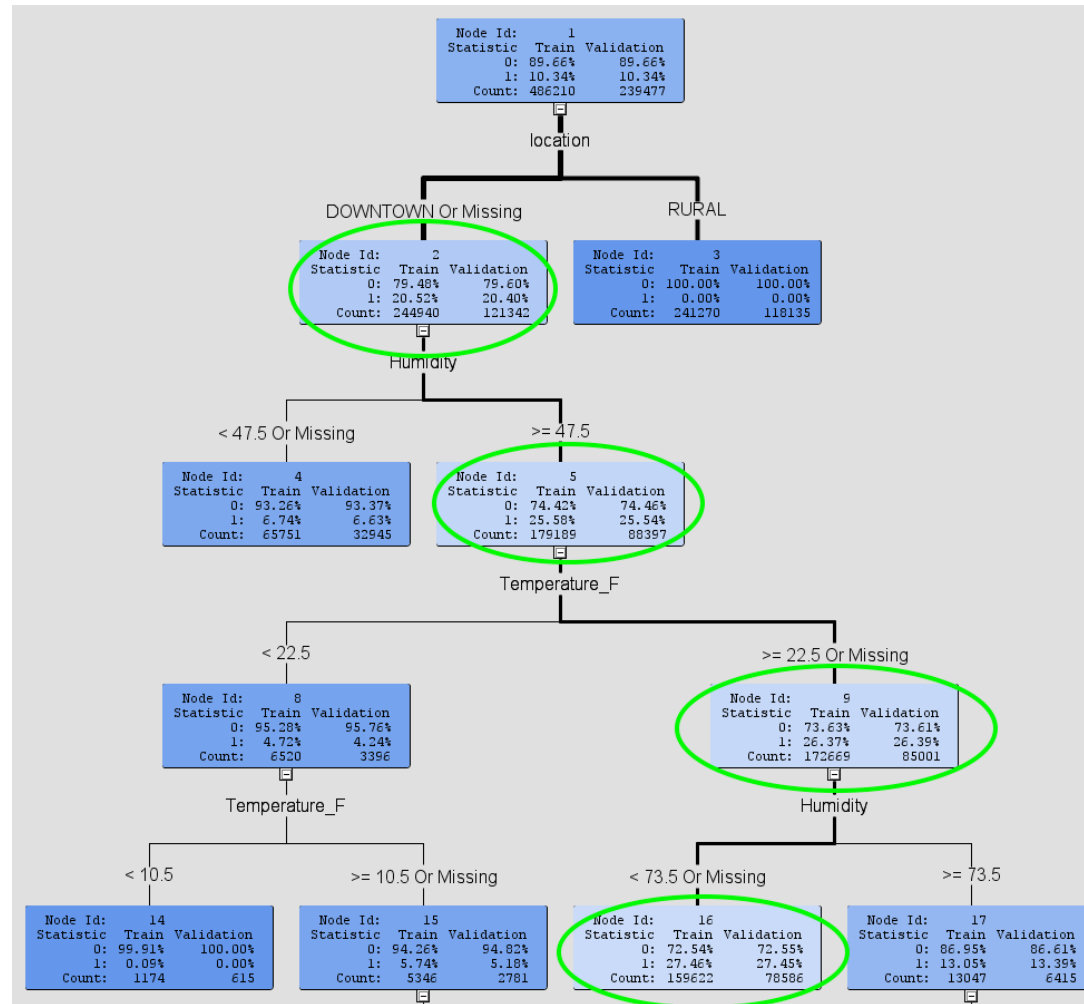


- SAS EM decided **Decision Tree 2** is best mathematically. Barely.
- Is this model also the most helpful when answering our question
 - **What affects air quality?**
 - **Can we predict if air quality is safe using only data readily available on any weather report?**
- *Yes!* The decision tree creates a nice way to see if the air quality is Safe or Unsafe.

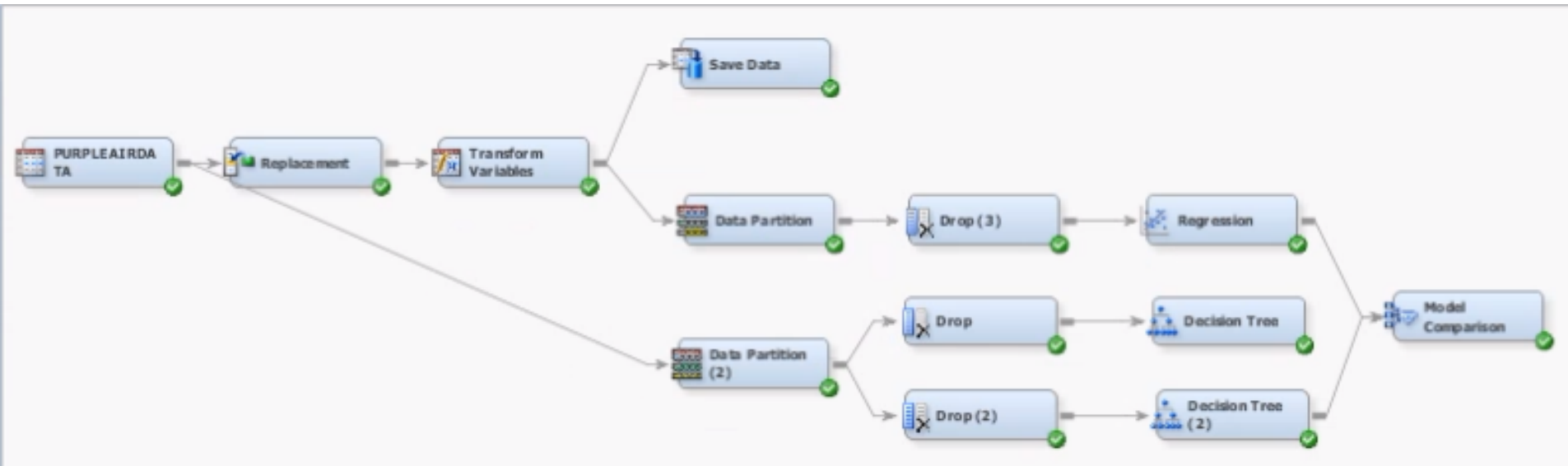


Modeling Method: Model Comparison

- Example:
It is 65 degrees Fahrenheit, 50% humidity, and you live 5 blocks from downtown.
- 27% Safe
73% Unsafe
- Consider YOUR context to make a decision!



+ Final SAS EM Diagram



+ Sources of Bias



- **Momentary readings** vs. EPA AQI is 24-hour average of readings
 - Biased towards spikes, which can be good or bad
- **Low cost sensor**
 - Not as sophisticated, finessed, or calibrated as other sensors
 - South Coast Air Quality Management District:
read high, specifically PM 2.5 levels overestimated by 36-48%
- **Local conditions at sensor site**
- **Measure of PM only!** The five major air pollutants regulated by the Clean Air Act are ground-level ozone, **particulate matter**, carbon monoxide, sulfur dioxide, and nitrogen dioxide.
Not a holistic picture.

+ Results & Conclusions

■ Temperature and Humidity do matter

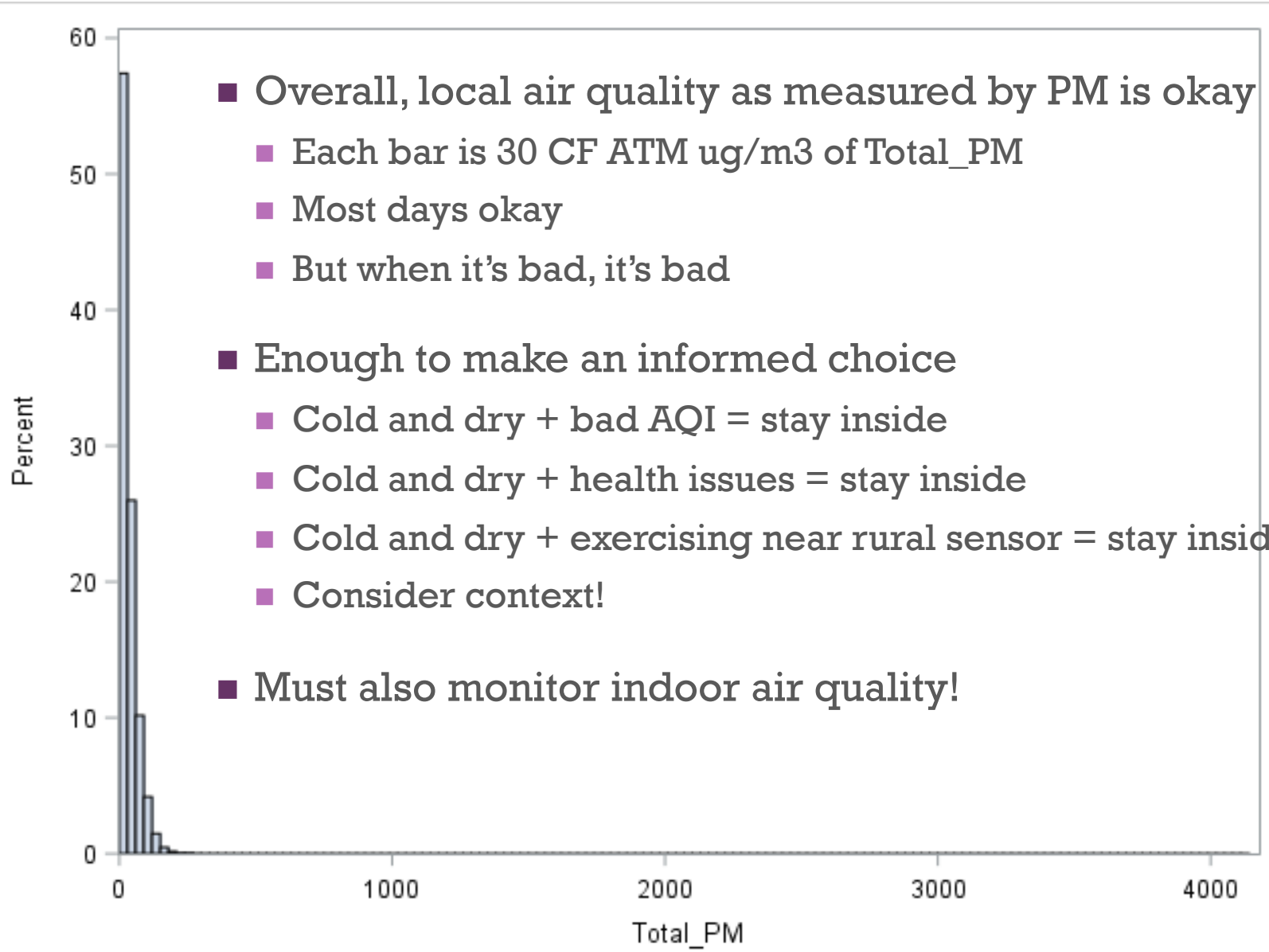
- Colder = More Likely to be Unsafe, Warm = More likely to be Safe
- Dry = More Likely to be Unsafe, Humid = More Likely to be Safe
- Cold and Dry is BAD! Warm and Humid is BETTER (but not great)!

■ Location matters a lot

- Less downtown/rural and more the DETAILS of that location
- Cause spikes that linger
- Decorah, IA
- Decorah 2



+ Results & Conclusions



- Overall, local air quality as measured by PM is okay
 - Each bar is 30 CF ATM ug/m3 of Total_PM
 - Most days okay
 - But when it's bad, it's bad
- Enough to make an informed choice
 - Cold and dry + bad AQI = stay inside
 - Cold and dry + health issues = stay inside
 - Cold and dry + exercising near rural sensor = stay inside
 - Consider context!
- Must also monitor indoor air quality!

+ Recommendations & Future Study



- Make informed choices! Do what you're comfortable with!
 - Tradeoffs!
- Need to study COMPLETE picture of Air Quality
 - Purple Air sensors monitor only one element of AQI
 - Need affordable sensors that monitor and report ALL 5 elements
 - Individuals have different sensitivities
 - Re-do analysis with more AQI factors included, more data, different locations
 - Might find different results
- *A decent AQI doesn't mean we shouldn't work towards cleaner air!*



+ Thank You!

