



FINAL REPORT  
GTI PROJECT NUMBER 22509-3

---

## 2019 Emission Factor Pilot Study

**Date Submitted**  
August 21, 2020

**Prepared For**  
Operations Technology Development (OTD) - SoCalGas

Ed Newton  
Research & Materials Manager  
Aviation Services  
Gas Engineering  
Southern California Gas Company  
San Diego Gas & Electric  
8101 S. Rosemead Blvd  
Pico Rivera, CA 90660  
ML: SC722S  
Tel:213.244.4238; Cell:213.219.0373  
ENewton@socalgas.com

**Project Manager**  
Kristine Wiley  
GTI  
847.768.0910  
kwiley@gti.energy

**Principal Investigator / Primary Author**  
Daniel Ersoy  
Element Resources, LLC  
847.343.9755  
dersoy@elementresourcesllc.com

**Project Team Members**  
Daniel Ersoy, Ed Newton, Jerone Powell, Gerry Bong and Kristine Wiley

GTI  
1700 S. Mount Prospect Rd.  
Des Plaines, Illinois 60018  
www.gti.energy



---

OTD Approved for Public Release

# Legal Notice

This information was prepared by Gas Technology Institute ("GTI") for Operations Technology Development/SoCalGas.

Neither GTI, the members of GTI, the Sponsor(s), nor any person acting on behalf of any of them:

Makes any warranty or representation, express or implied with respect to the accuracy, completeness, or usefulness of the information contained in this report, or that the use of any information, apparatus, method, or process disclosed in this report may not infringe privately-owned rights. Inasmuch as this project is experimental in nature, the technical information, results, or conclusions cannot be predicted. Conclusions and analysis of results by GTI represent GTI's opinion based on inferences from measurements and empirical relationships, which inferences and assumptions are not infallible, and with respect to which competent specialists may differ.

Assumes any liability with respect to the use of, or for any and all damages resulting from the use of, any information, apparatus, method, or process disclosed in this report; any other use of, or reliance on, this report by any third party is at the third party's sole risk.

The results within this report relate only to the items tested/reviewed.

# Table of Contents

	Page
<b>Legal Notice .....</b>	<b>i</b>
<b>Table of Contents .....</b>	<b>ii</b>
<b>List of Figures .....</b>	<b>vi</b>
<b>List of Tables .....</b>	<b>viii</b>
<b>Executive Summary .....</b>	<b>1</b>
<b>1. Introduction - Report Layout.....</b>	<b>7</b>
<b>2. Background - Regulator, Industry Studies, and SoCalGas Process Development .....</b>	<b>9</b>
2.1. California Rule Making .....	9
2.2. Leak Grading in California .....	9
2.3. Summary of Studies Referenced or Developed as Part of this Report .....	11
Data Obtained from Past Industry Studies as Comparisons for this Report .....	12
SoCalGas Studies for this Report – Data and Study Terminology Definitions.....	12
<b>3. Approach for Field Sampling, Measurement Techniques, and Decision Tree Process .....</b>	<b>14</b>
3.1. SoCalGas Process Development.....	14
Initial focus on Existing System Data and Leak Centering Process.....	14
Enhanced Leak Survey Practice .....	15
3.2. Development of the Decision Tree Approach .....	15
Types of Surface Conditions .....	15
Methane Concentration Measurement Process .....	16
Decision Tree Concentration Threshold Values for Leak Flow Rate Measurement .....	17
Leak Flow Rate Measurement Process .....	17
3.3. Selection of Single vs. Facility/Material Specific EFs .....	18
3.4. Precision and Sample Size Analysis - Minimum Sample Size .....	18
Sample Size for a One Mean Confidence Interval.....	18
Sample Size for a Bayesian Proportional Analysis.....	20
3.5. Physical Sampling Plan - Where to Sample.....	21
<b>4. Methodology Overview of Data Collection and Statistical / Probabilistic Analysis .....</b>	<b>22</b>
4.1. Descriptive Statistics of Study Samples.....	22
4.2. Data Transformation, Regression, and MCMC Models.....	22
4.3. Decision Tree Predictive Performance.....	22
4.4. Population Mean Leak Rate Inferential Analysis .....	22
4.5. Emission Factor Determination.....	22

4.6. Method of Emission Factor Application .....	23
4.7. Quality Assurance .....	23
Statistical Checks .....	23
Probabilistic Regression Check .....	23
Significant Figures .....	23
Standard Conditions .....	24
Distribution Fitting and Monte Carlo Sampling .....	24
Leak Concentration and Flow Rate Measurement Error/Uncertainties .....	24
Carrying Uncertainty Through to the Emission Factor Calculations .....	25
<b>5. Analysis and Results.....</b>	<b>26</b>
5.1. Descriptive Statistics.....	26
Comparison of SoCalGas Study to Industry Studies .....	26
Data analysis using Dot Plots .....	28
Data analysis using Box Plots .....	31
5.2. Removal of Studies with Sample Bias .....	33
Data Analysis using Cumulative Fraction Plots.....	35
5.3. Data Transformation .....	38
Histogram of Log Transformed Leak Data .....	38
One-sample Kolmogorov-Smirnov Test for Lognormality .....	39
Quantile-Normal Plot Analysis of Transformations .....	40
5.4. Linear Regression Study Means Comparison .....	42
Combined National vs. Combined SoCalGas Study Means Analysis .....	42
Analysis of Variance (ANOVA).....	42
Linear Regression .....	44
Overview .....	44
Three SoCalGas Study Means Analysis .....	45
Analysis of Variance (ANOVA).....	45
Linear Regression .....	46
Linear Regression Residual Analysis and Regression Diagnostics .....	47
5.5. Bayesian Monte Carlo Markov Chain (MCMC) Regression .....	47
5.6. Sensitivity of Leak Rate to Geographic District and Year of Detection Analysis .....	48
Geographic District of Leak.....	48
Year of Leak Detection .....	48
5.7. Concentration vs. Leak Rate Analysis .....	53
General Trends.....	53

5.8. Decision Tree Leak Prediction Quantitative Performance .....	55
Leak Rate Statistics of Empirical Data by Decision Tree Groupings .....	55
Overall Results from Empirical Data.....	55
Leak Rate Statistics of Empirical Data by Confirmed Leak Rate Groupings .....	56
Actual Leak Rate less than 10 scfh and True and False Negatives.....	56
Actual Leak Rate greater than or equal to 10 scfh .....	58
5.9. Bayesian Probabilistic Decision Tree Error-Type Analysis .....	60
Joint False/True Positive (Type I) and Negative (Type II) Errors.....	60
Independent False/True Positive Error Type I.....	60
Independent False/True Negative Error Type II.....	61
Overall Efficiency of Decision Tree Process.....	62
2019 3-District Pilot Data Example .....	62
Additional Ongoing False Negative Validation Sampling.....	63
5.10. Population Mean Leak Rate Analysis .....	64
Bootstrap Analysis of Field Leak Rate Data.....	64
Monte Carlo Analysis of Fitted Distribution (for illustrative purposes only) .....	64
Mean Leak Rate Analysis Results .....	64
Log-Normal Distribution Fit and Monte Carlo Analysis.....	66
<b>6. Emission Factor Development and Application .....</b>	<b>68</b>
6.1 Development of a Company Specific Emission Factor.....	68
Input Information .....	68
Calculation of Distinct Emission Factors .....	69
6.2. Table of Emission Factors .....	69
6.3. Carrying Uncertainty Through to the Emission Factor Calculations .....	71
6.4. Scenarios of EF Application .....	72
<b>7. Summary of Results and Conclusions .....</b>	<b>73</b>
<b>Appendix A: Surface Measurements of Underground Leak Flow Rate.....</b>	<b>76</b>
<b>Appendix B: Statistical and Probabilistic Analysis Details and Supplemental Analysis .....</b>	<b>78</b>
Linear Regression Residual Analysis and Regression Diagnostics .....	78
Residual Distribution .....	80
Lowess .....	80
DFBETA.....	81
Cook's Distance.....	81
Leverage .....	82
Bayesian Monte Carlo Markov Chain (MCMC) Models of Sample Leak Rates.....	82
Metropolis-Hastings Sampling (MHS) .....	83

Gibbs Sampling (GS) .....	84
ANOVA and Pairwise comparison of National and SoCalGas Studies .....	86
Linear and Logistic Regression of Concentration vs. Leak Flow Rates .....	87
Linear Regressions .....	87
Logistic Regressions .....	91
<b>Appendix C: Log-normal Distribution Facts .....</b>	<b>97</b>
Log-normal Distribution Equations.....	97
Log-normal Plot and Comparison to Normal Distribution.....	97
Distribution Fitting Goodness of Fit .....	98
<b>Appendix D: Leak Spread Comparison to Leak Rate .....</b>	<b>99</b>
<b>Appendix E: Study Leak Rate and Concentration Data .....</b>	<b>102</b>
<b>References .....</b>	<b>114</b>

# List of Figures

	Page
<b>Figure 1: Timeline of Major Tasks and Milestones.</b> .....	12
<b>Figure 2: Surface Condition Location Definitions.</b> .....	16
<b>Figure 3: Sample Size for a One Mean Confidence Interval of Log(10) Leak Rate.</b> .....	19
<b>Figure 4: Sample Size for a Bayesian Proportional Analysis by Confidence Level.</b> .....	20
<b>Figure 5: Leak Rate Median and Mean Plot by National and SoCalGas Study.</b> .....	27
<b>Figure 6: Leak Rate Plot by Sample for Five SoCalGas Studies.</b> .....	28
<b>Figure 7: Leak Rate Plot by Sample for Three National Studies.</b> .....	29
<b>Figure 8: Leak Rate Lower/Upper Percentiles by Combined National and SoCalGas Studies.</b> ....	31
<b>Figure 9: Leak Rate Box Plots by National and SoCalGas Studies.</b> .....	32
<b>Figure 10: Leak Rate Lower and Upper Percentiles by SoCalGas Study.</b> .....	33
<b>Figure 11: Leak Rate Cumulative Fraction of Two Sample-Biased SoCalGas Studies.</b> .....	34
<b>Figure 12: Leak Rate Cumulative Fraction of Combined National and SoCalGas Studies.</b> .....	35
<b>Figure 13: Leak Rate Cumulative Fraction of National Studies.</b> .....	36
<b>Figure 14: Leak Rate Cumulative Fraction of Three SoCalGas Studies.</b> .....	37
<b>Figure 15: Leak Rate Histogram of Log(10) of Combined National and SoCalGas Studies.</b> .....	39
<b>Figure 16: Quantile-Normal Plot of Log-normal Transformed Combined National Studies.</b> ....	41
<b>Figure 17: Quantile-Normal Plot of Log-normal Transformed Combined SoCalGas Studies.</b> ...41	41
<b>Figure 18: Leak Rate Box Plots by Year Leak Detected for Three Combined SoCalGas Studies.</b> 49	49
<b>Figure 19: Leak Rate Median, Mean, and Maximum by Year SoCalGas Leaks Detected.</b> .....	50
<b>Figure 20: Separate Plots of Leak Concentrations vs. Rates by DT Category.</b> .....	53
<b>Figure 21: Mean Leak Rate for DT Categories.</b> .....	56
<b>Figure 22: Mean Leak Rate for DT Categories when Actually &lt; 10 scfh.</b> .....	57
<b>Figure 23: Mean Leak Rate for DT Categories when Actually <math>\geq</math> 10 scfh.</b> .....	58
<b>Figure 24: Expected Decision Tree Output with No Concentration Data.</b> .....	61
<b>Figure 25: Expected Decision Tree Output with Known DT Category.</b> .....	62
<b>Figure 26: Leak Rate Cumulative Fractions of Combined SoCalGas Studies and L-N Fit.</b> .....	66
<b>Figure 27: Bootstrap Mean Leak Rate Box Plots of Field Data and Log-Normal Fit.</b> .....	67
<b>Figure 28: Quantifying surface flux rate of an underground emission.</b> .....	76
<b>Figure 29: Schematic of surface chamber measurements with the Hi-Flow sampler.</b> .....	77
<b>Figure 30: Histogram Diagnostic Plot of LR Leak Rate Residuals for Three SoCalGas Studies.</b> ..80	80
<b>Figure 31: Lowess LR Diagnostic of Leak Rate Residuals by SoCalGas Study.</b> .....	80
<b>Figure 32: DFBETA LR Diagnostic of Leak Rate for Three SoCalGas Studies.</b> .....	81

**Figure 33: Cook's Distance LR Diagnostic of Leak Rate for Three SoCalGas Studies. ....81**

**Figure 34: Leverage LR Diagnostic of Leak Rate for Three SoCalGas Studies.....82**

**Figure 35: Diagnostics for Bayesian MCMC(MHS) of Leak Rate Means - SoCalGas Studies.....85**

**Figure 36: Diagnostics for Bayesian MCMC(GS) of Leak Rate Means - SoCalGas Studies. ....85**

**Figure 37: Predictive Margins for Leak Rate by Concentration from Linear Model. ....89**

**Figure 38: Sensitivity of SSS Concentration to Large Leak Detection (Logistic Regression).....91**

**Figure 39: Sensitivity of CIP Concentration to Large Leak Detection (Logistic Regression).....92**

**Figure 40: Sensitivity of US Concentration to Large Leak Detection (Logistic Regression). ....92**

**Figure 41: Sensitivity of BH Concentration to Large Leak Detection (Logistic Regression).....93**

**Figure 42: Sensitivity of SSS Concentration to Large Leak Detection - Multi-Thresholds.....94**

**Figure 43: Sensitivity of CIP Concentration to Large Leak Detection - Multi-Thresholds.....95**

**Figure 44: Sensitivity of US Concentration to Large Leak Detection - Multi-Thresholds. ....95**

**Figure 45: Log-normal Distribution Theoretical Plot. ....97**

**Figure 46: Scatter Plot of Leak Concentration Spread vs. Leak Rate. ....99**

**Figure 47: Histogram of Leak Rates by Leak Concentration Spread..... 100**

**Figure 48: Box Plot of Leak Concentration Spread vs. Leak Rate. .... 100**

**Figure 49: Mean and Median of Leak Rate by Leak Concentration Spread..... 101**

# List of Tables

	Page
<b>Table 1: Summary of Studies Referenced and Performed as Part of this Report.....</b>	<b>11</b>
<b>Table 2: Sample Size for a One Mean Confidence Interval of Log(10) Leak Rate.....</b>	<b>19</b>
<b>Table 3: Leak Rate Mean, Minimum, and Maximum by National and SoCalGas Study. ....</b>	<b>27</b>
<b>Table 4: Leak Rate Mean and 95% C.I. National vs. SoCalGas Studies. ....</b>	<b>29</b>
<b>Table 5: Leak Rate Median, 5%, and 95% Percentiles by National and SoCalGas Study.....</b>	<b>30</b>
<b>Table 6: Kolmogorov-Smirnov test of Combined National Studies for Normality.....</b>	<b>40</b>
<b>Table 7: Kolmogorov-Smirnov test of Combined SoCalGas Studies for Normality. ....</b>	<b>40</b>
<b>Table 8: ANOVA of Combined National vs. Combined SoCalGas Leak Rate Means.....</b>	<b>43</b>
<b>Table 9: Linear Regression of Individual National and SoCalGas Study Leak Rate Means. ....</b>	<b>44</b>
<b>Table 10: ANOVA of Leak Rate Means for Three SoCalGas Studies. ....</b>	<b>45</b>
<b>Table 11: LR and PW Comparison of Leak Rate Means for Three SoCalGas Studies. ....</b>	<b>46</b>
<b>Table 12: ANOVA of Leak Rate Means Across Districts for Three SoCalGas Studies. ....</b>	<b>48</b>
<b>Table 13: ANOVA of Leak Rate Means Across Year Detected for Three SoCalGas Studies. ....</b>	<b>50</b>
<b>Table 14: PW Comparison of Leak Rate Means by Year Detected for Three SoCalGas Studies. .</b>	<b>52</b>
<b>Table 15: Combined Bootstrap Mean and C.I. for Three SoCalGas Studies (Baseline).....</b>	<b>55</b>
<b>Table 16: Leak Rate Mean, Min., and Max. of by DT Grouping. ....</b>	<b>55</b>
<b>Table 17: Leak Rate Median, 5%, and 95% Percentiles by DT Grouping. ....</b>	<b>56</b>
<b>Table 18: Leak Rate Mean, Min., and Max. of Confirmed &lt;10 scfh by DT Grouping.....</b>	<b>57</b>
<b>Table 19: Leak Rate Med., 5%, and 95% Percentiles of Confirmed &lt;10 scfh by DT Grouping....</b>	<b>57</b>
<b>Table 20: Bootstrap Leak Rate Mean and C.I. for Confirmed &lt; 10 scfh (Actual Negatives). ....</b>	<b>57</b>
<b>Table 21: Leak Rate Mean, Min., and Max. of Confirmed ≥10 scfh by DT Grouping.....</b>	<b>58</b>
<b>Table 22: Leak Rate Med., 5%, and 95% Percentiles of Confirmed ≥10 scfh by DT Grouping....</b>	<b>58</b>
<b>Table 23: Bootstrap Leak Rate Mean and C.I. for Confirmed ≥ 10 scfh (Actual Positives).....</b>	<b>59</b>
<b>Table 24: Type I &amp; II Uninformed (DT Cat. Unknown) Errors with 5% and 95% Pred. Limits. .</b>	<b>60</b>
<b>Table 25: Type I Errors with 5% and 95% Prediction Limits for DT Positive Group. ....</b>	<b>60</b>
<b>Table 26: Type II Errors with 5% and 95% Prediction Limits for DT Negative Group. ....</b>	<b>61</b>
<b>Table 27: Leak Rate Bootstrap Means by Study Group.....</b>	<b>65</b>
<b>Table 28: Leak Rate (scfh) Percentiles of the Bootstrap Mean.....</b>	<b>65</b>
<b>Table 29: Table of SoCalGas Company Specific Emission Factors by DT Grouping. ....</b>	<b>69</b>
<b>Table 30: Table of Emission Factors to use for Field Situations.....</b>	<b>72</b>
<b>Table 31: Bayesian MCMC(MHS) of Leak Rate Means - Three SoCalGas Studies. ....</b>	<b>83</b>
<b>Table 32: Bayesian MCMC(GS) of Leak Rate Means for Three SoCalGas Studies. ....</b>	<b>84</b>
<b>Table 33: ANOVA of Individual National and SoCalGas Study Leak Rate Means.....</b>	<b>86</b>

**Table 34: PW Comparison by National/SoCalGas Study Pair for Leak Rate Means. ....87**  
**Table 35: Linear Regression of DT Concentration and SoCalGas Study Leak Rates.....90**  
**Table 36: Log-Normal Distribution Equations.....97**  
**Table 37: Distribution Goodness of Fit to SoCalGas data set (291 samples). ....98**  
**Table 38: Leak Spread Categories. ....99**  
**Table 39: Leak Rate scfh Methane (CH<sub>4</sub>) and Concentration (% gas) by Study..... 102**

# Executive Summary

## Objective

Traditionally, inventories and reporting programs use emission factors based on leaks/emissions expected *per mile* of a specific *type* of pipe. This means that currently, the only way to show reductions is to reduce the number of miles of higher emitting types of pipe, such as cast iron and cathodically unprotected steel pipe. More recent regulations have included emission factors based solely on leaks.

The objective of this study was to develop a method for flagging large leaks for cost-effective measurement and repair to minimize systemwide methane leakage rates that further focuses on non-hazardous (grade 2 and 3) leaks. So, if a company can reduce its number of higher emitting non-hazardous leaks, it can reduce actual emissions *and* more accurately demonstrate the reduction.

The study conducted statistically sound sampling of non-hazardous pipeline leaks using well-proven field measurement techniques to provide data to calculate company-specific methane emission factors for Southern California Gas Company (SoCalGas) buried Distribution system leaks. State-of-the-art parametric and non-parametric statistical analysis, resampling, Monte Carlo, and Bayesian probabilistic analysis were used when appropriate.

The approach demonstrates that methane concentrations collected at designated types of locations, in the manner prescribed, and analyzed according to the Decision Tree process can be used to predict whether a leak flow rate is either above or below a specified target flow rate.

## Background - Regulator, Industry Studies, and SoCalGas Process Development

Based on 2015 California state rulemaking, the research team at SoCalGas began the development of an approach for identifying and differentiating leaks on the buried distribution system that have relatively high flow rates, for the purpose of prioritizing repairs and reducing natural gas emissions from the distribution system. The approach was chosen with the objective of developing a “cost-effective” methodology as defined within the California State CPUC Rulemaking (R.) 15-01-008 [1]. Various prior industry studies, which include leak data from the SoCalGas system, were leveraged for this effort.

Initially, the work was based on the hypothesis that a small percentage (approximately 5%) of non-hazardous buried leaks in the SoCalGas distribution system had a flow rate of 10 scfh or larger, and that existing system data about the leaks obtained at the time the leaks are detected and graded could be leveraged to identify a sub-set of all system leaks that had the greatest probability of being high flow-rate leaks. Since leaks that are categorized as a safety hazard (“Code 1” or “Grade 1” leaks) are identified more readily and fixed immediately, studies have focused on the non-hazardous leaks (“Code 2 or 3” leaks) that are generally scheduled for later repair or monitored, and thus can continue emitting for a longer period of time. Prior industry studies suggested that the number of system non-hazardous leaks that have a high flow rate is a

small percentage of total leaks. The goal was to verify the flow rate population distribution of non-hazardous leaks in the SoCalGas system and find an efficient way to identify large flow rate non-hazardous leaks, so they can be prioritized for repair.

SoCalGas mined existing system data but found that there was no discernable relationship between available data and the flow rate of the leak. After much research, it was determined that surface expression measurements across the entire spread of the leak were needed to determine the flow rate. As more surface expression leak flow rate data was collected, some relationships to methane concentration measurements began to emerge.

Extensive surface expression method leak flow rate data was then collected along with methane concentration data at the corresponding prescribed surface measurement locations. Subsequent groupings of that data based on similar location descriptions yielded promising correlations to the leak flow measurements.

Starting in 2016, surface expression measurements were analyzed against site-specific characteristics of ground-level methane concentration measurements, and concentration thresholds were developed for each surface category to identify the leaks with the potential for a high flow rate. A process and workflow were then developed where surface expression measurements to calculate leak flow rate are performed whenever one or more of the concentration threshold values are met or exceeded, so repairs can be prioritized. This methodology was termed the “Decision Tree” (DT) approach.

In order to validate and achieve high statistical confidence in the DT model output a statistical and probabilistic data analysis study commenced, with the results presented in this report. This work led to the collection of leak flow data based on a geographically diverse random sample of the entire SoCalGas distribution system which established a technically sound foundational leak flow rate dataset. The data set used for the emission factor calculations in this report included 291 such samples.

## **Approach - Field Sampling, Measurement Techniques, and Decision Tree**

### *System-Wide Random Field Sample (Leak Site) Design*

For this pilot effort, SoCalGas stratified its sample population by district, and then randomly drew the corresponding number of leak site samples (per district or district grouping) to preserve the correct proportion of the districts in the total population of leaks.

### *Concentration Measurements*

The operator utilizes either the Heath DP-IR (Detecto Pak Infrared) or GMI Gasurveyor along with the survey probe attachment to survey the leak site. The spread of the leak is determined by probing the ground surface and identifying the extent at which any methane concentration is present. Once the spread is determined the operator then identifies and records the highest sustained ground-level reading within the spread of the leak for each of the four surface conditions where gas indications are found.

### *Decision Tree Process*

A process was then developed where surface expression measurements to calculate leak flow rate are performed whenever one or more of the leak concentration threshold values are met or exceeded, so repairs can be prioritized.

This methodology was termed the “Decision Tree” (DT) approach. The applicable concentration measurements are compared to the threshold values, and any one (or more) of the concentration measurements taken from the four prescribed surface conditions, that meets or exceeds the threshold concentration value will then result in that leak being classified as a potential large, non-hazardous leak with a possible leak flow rate of 10 scfh or higher.

The threshold values for concentration measurements by surface condition type are:

- 20% Gas: Crack (or seam) In Pavement - CIP
- 5% Gas: Unpaved Surface - US
- 80% Gas: Bar Hole (leak survey type) - BH
- 60% Gas: Small Sub-Structure (not gas system related) - SSS

### *Leak/Emission Flow Rate Measurements*

The leakage flow rates were measured using the well-established and published surface expression methodology. These leak rate measurements provide an approximation of ‘in-air’ methane emission rates without the need to excavate the leak source.

## **Methodology - Data Collection and Statistical/Probabilistic Analysis**

Non-hazardous leak survey data, including methane concentration measurements at defined types of surface condition locations, were collected as part of this study. The specific numbers of the various sampling efforts are listed in the Background section. Leakage flow rates were measured from selected underground distribution pipeline leaks, triggered based on the Decision Tree concentration thresholds, and based on a random sample across the entire SoCalGas distribution system.

Standard descriptive statistical analysis was conducted including calculation of sample means, medians, percentiles, inner quartile ranges, and other statistics. Various analysis and plotting techniques were used to confirm sampling bias and draw high-level conclusions on the different individual and grouped sample leak rate distribution center tendencies, uncertainties, and shape.

Data transforms were used to ensure that any regression model utilized had a sound basis. Monte Carlo Markov Chain (MCMC) with Metropolis-Hastings Sampling (MHS) and Gibbs Sampling (GS), Linear Regression (LR), and Analysis of Variance (ANOVA) were used to quality check sample set data, spot outliers, confirm assumptions, assess regression, and check probabilistic residuals and diagnostic measures.

A purely probabilistic Bayesian analysis was used to measure the Decision Tree performance by grouping leaks into two categories. By using this approach and analysis, there is no model “form” that needed to be “informed” or “trained”. The analysis incorporated and related the leak

concentration levels with the Decision Tree threshold point at 10 scfh between “Large” and “Not Large” leak groups. The DT performance metrics included developing a False/True Negative/Positive (Type I and II) error table.

Resampling with replacement (bootstrap) analysis of field leak rate data and Monte Carlo sampling of a fitted data distribution of leak rates were both used to infer the population mean leak rates with upper and lower confidence limits from the sample data. The SoCalGas emission factors were derived using a combination of the appropriate bootstrap population leak rate means and the Bayesian Decision Tree error table percentiles.

## **Summary of Results, Emission Factor Development, and Application**

The national studies compared well with the SoCalGas studies. The mean, median, and upper and lower 95% percentiles for leak rate of these two groups are similar.

Two of five SoCalGas sample sets were known to contain sample bias, as well as being an order of magnitude in size smaller than the other three. These were analyzed in this report to show how bias might appear during analysis, and they were not included in the ultimate combined data set.

The non-hazardous leak rate values from the SoCalGas combined data set was analyzed for unexplainable outliers or extreme values and was log transformed, resulting in a normally distributed data set. Upon review of the extreme values, all of them were deemed as sound data points and not errors or anomalous values. The log-normal transformation of the leak rate data permitted a variety of statistical regression tools to be appropriately leveraged.

A series of regression and probabilistic analysis were conducted on the data set. Two key findings were that when the samples sizes supported categorical analysis that there was no significant sensitivity of the leak rate means to geographic operating districts where the leak was found, or the time interval from when the leak was detected.

An analysis of the field methane concentration vs. measured leak rates was done by Decision Tree methane concentration threshold category. The regression analysis of the mean leak flow rate vs. methane concentration showed the expected upward trend for the average values. The concentration threshold intersection with the established 10 scfh “Large” vs. “Not Large” flow rates were within the 95% confidence interval of the regression model or above and to the left (a conservative situation) of the predictive margin plots.

A Bayesian probabilistic analysis was conducted of the Decision Tree threshold performance. This resulted in a true/false-positive/negative Error Table. The Decision Tree thresholds correctly assigned low leak situations 98.9% of the time, i.e. true negatives with a 95% prediction interval of 98.9% to 99.5%. Likewise, the Decision Tree had a false negative (Type II error) of only 1.1% with a 95% prediction interval of 0.47% to 3.6%.

The leak rate data was bootstrapped 10,000 times with replacement and a re-sample size equal to the field data sample size. This analysis provided the overall mean leak rate, as well as the mean

leak rates for less than ( $<$ ) 10 scfh leakers and greater than or equal to ( $\geq$ ) 10 scfh leakers - all from the empirical data. The bootstrap analysis provided a full set of percentiles for the actual mean leak rates which allows one to establish confidence intervals for the mean values at any desired confidence level.

The leak rate data was fit to a log-normal distribution as well, and this fit was used to conduct a Monte Carlo analysis of the mean leak rates as was conducted with the bootstrap analysis using the actual field leak rate data. The same re-sample and over sample sizes were used as was done with the bootstrap analysis to properly propagate the uncertainty through the analysis. The result showed the two approaches were very similar, with the Monte Carlo of the log-normal distribution fit being conservative in the low- to mid- leak rate ranges and about the same in the high range of leak rates.

A set of emission factors based on the Decision Tree categorization were calculated by combining the mean leak rates with their corresponding expected percentiles (in a weighted manner) from the Decision Tree error table. It was noted that the Decision Tree derived emission factors were conservative (higher) than one would have obtained from a straight average of the empirical data from the All District Study of the SoCalGas system. This is due to the Bayesian analysis properly accounting for false negatives in the Decision Tree process.

A calculation of the efficiency of the process was done using the 2019 3-District Pilot study which had a total number of 356 screened leaks with surface concentration measurements. Of these, the DT was triggered for flow rate measurement 44 times. This therefore relates to a flow rate measurement ratio of  $44 / 356$  or 12.4%, meaning that when considering leak sites visited and screened with surface concentration measurements that one would expect to be triggered by the DT process and criteria to have approximately 1 in 8 of them classified as potential non-hazardous large leak rates and be scheduled for leak rate measurement or prioritized for repair.

For this particular example, rather than measuring all 356 leaks to find all the large leaks; the DT process was used resulting in the requirement to measure only 1 in 8 leaks while maintaining a false negative error of 1.1%. In summary:

- Using the DT method, 4 of the expected 7 large leaks were found by measuring the leak flow rate from 44 out of 356 leak sites.
- Without the DT, to find the same ratio of 4 out of the 7 large leaks, 203 leak flow rates on average would need to be measured out of the 356 leak sites.
- This means the DT efficiency increase is  $203/44 = 4.6x$  (460%) more efficient at finding the same number of large leaks when not using the DT process.
- The DT is therefore an efficient screening mechanism, with a high potential to continue to improve over the short-term full implementation period.

## Conclusions

SoCalGas conducted a statistically sound study of pipeline leaks using random samples and well-proven field leak concentration and flow rate measurement techniques to provide data to

calculate SoCalGas company-specific natural gas emission factors for buried distribution system non-hazardous leaks.

The developed Decision Tree approach of using concentration measurements with thresholds to establish large and not large non-hazardous leaks was successful as measured by a 98.9% true negative value associated with predicted leak and actual leak rates.

The inferred population mean leak rates were combined with the associated Decision Tree performance percentages to calculate appropriately weighted emission factors for large and not large non-hazardous leaks.

This allows the assignment of emission factors for the not large non-hazardous leaks that would not have leak rate flow measurements performed on them, as well as any Decision Tree classified large non-hazardous leaks that did not have leak rate flow measurements performed.

The approach will be further refined and improved by continuing to:

- Collect field data leading to lower uncertainty, i.e. tighter confidence intervals around leak and Decision Tree performance metrics;
- Perform random checks for false negatives to identify possible upset conditions in expected leak rates, e.g. from a change in system performance and/or environmental stressors; and
- Analyze and adjust the Decision Tree thresholds or even add new thresholds to further increase the method's predictive accuracy and/or increase process efficiency to continuously improve the cost-effectiveness of the approach, overall process for detection, and repair of large flow system leaks to minimize natural gas emissions.

# 1. Introduction - Report Layout

---

Below is a brief description of the major sections of this report.

- 1. Introduction - Report Layout.** This section describes the major sections of the report and their content.
- 2. Background - Regulator, Industry Studies, and SoCalGas Process Development.** This section discusses national and California rulemaking related to natural gas emissions. A table of the past industry emission studies and the data sets used as part of this study is provided, along with details on the current studies, sample sets, and emission factor calculations used in this report. A basic timeline is also provided showing the order of activities in developing the Decision Tree process in relation to the various sets of data.
- 3. Approach for Field Sampling, Measurement Techniques, and Decision Tree Process.** This section contains the development progression and components of the Decision Tree process (method). It includes how and when surface concentration measurements are taken, how leak rate size is initially estimated, and how leak flow rates are measured. A precision and sample size analysis related to the desired confidence interval width for mean leak rate is presented, as well as a minimum sample size for Bayesian probabilistic analysis used to measure the performance of the Decision Tree method. A discussion on random sampling is also presented.
- 4. Methodology Overview of Data Collection and Statistical / Probabilistic Analysis.** The analysis methods employed in the Analysis and Results section are listed, and details are provided. A quality assurance section lists the statistical checks, secondary analysis, significant figure management, and standard conditions of the data collection. The methods for distribution fitting and Monte Carlo sampling of the same are described. Details on the leak concentration and leak flow rate measurement uncertainties, as well as the inferential statistical and probabilistic analysis uncertainties are addressed.
- 5. Analysis and Results.** This is the largest section of the report and documents the collection and analysis of the SoCalGas study data. The leak rate data from the SoCalGas data sets is analyzed and compared with the national industry studies listed in the Background section. Descriptive statistics are calculated, and data is plotted. The data is transformed and checked for normality in its distribution graphically and with non-parametric statistical tests. Analysis of variance and linear regression are used to compare the leak rates of various studies as well as look at the relationships of leak surface concentration values to leak rate values. The Decision Tree prediction performance is analyzed quantitatively with a Bayesian probabilistic analysis and the average leak rate and population distributions of the various studies are calculated through bootstrap resampling of the field data. The leak rate data sets

are also fit to a log-normal distribution to demonstrate the ability to "stretch" small data sets as stopgap measure until adequate sample sizes are achieved for bootstrap means analysis.

- 6. Emission Factor Development and Application.** In this section the bootstrapped leak rate distributions are combined with the Decision Tree prediction performance metrics to calculate company specific emission factors.
- 7. Summary of Results and Conclusions.** This section contains a high-level summary of the results and conclusions of the study and recommended next and ongoing steps.

**Appendix A: Surface Measurements of Underground Leak Flow Rate.** This section provides details on the well-established technique used in the study for surface measurements of underground leak flow rates.

**Appendix B: Statistical and Probabilistic Analysis Details and Supplemental Analysis.** This section lists the technical details, detailed output, diagnostics, and residual analysis of the Bayesian Monte Carlo Markov Chain (MCMC) models, select linear regression, and logistic regression analyses.

**Appendix C: Log-normal Distribution Facts.** The details on the log-normal distribution, as well as the goodness-of-fit measures for the fit of the SoCalGas study data set are listed in this section.

**Appendix D: Leak Spread Comparison to Leak Rate.** Additional details on attempts to correlate the spatial spread of leaks concentration measurements to leak rate are listed in this section.

**Appendix E: Study Leak Rate and Concentration Data.** The leak rate and concentration observations for the national and SoCalGas studies are listed in this section.

**References.** The references cited in the report are listed in this section.

## 2. Background - Regulator, Industry Studies, and SoCalGas Process Development

---

Since the Natural Gas industry first began studying distribution system leak rates in the early 1990's SoCalGas has been a leader in both funding the work as well as providing its resources and facilities to conduct and participate in the studies. Data obtained by independent researchers from SoCalGas system leaks have been used in the following industry studies: EDF/WSU – 2015[1], GTI/OTD – 2013[2], DOE/GTI – 2019[3], CARB/GTI – 2019[4].

### 2.1. California Rule Making

Historically, public safety has been the driver for California gas utilities policy and procedures for identifying and repairing distribution system leaks that are potentially hazardous as soon as reasonably possible. However, on January 22, 2015, the CPUC opened Rulemaking (R.) 15-01-008 [5] to implement the provisions of Senate Bill (SB) 1371 (Statutes 2014, Chapter 525) [6].

SB 1371 required the adoption of rules and procedures to minimize natural gas leakage from CPUC-regulated natural gas pipeline facilities as a means of reducing emissions of greenhouse gases. SB 1371 directs the Commission to consult with the California Air Resources Board (CARB) [7], to achieve the goals of the Rulemaking. California's statutory methane emissions reduction target is to lower 2030 levels to at least 40% below 2015 levels.

The SB 1371 Leakage abatement program uses Emission Factors for distribution mains and services from the 1996 GRI/EPA study, Methane Emissions from the Natural Gas Industry, GRI-94/0257.25, EPA-600/R-96-080, June 1996. Volume 9: Underground Pipelines[8]. These are leaker-based emission factors with engineering units of "Mscf NG/day/leak".

EPA Subpart W[9] uses GRI-GHGCalc[10] population emission factors to estimate emissions from distribution mains and services derived from the same report used for the SB 1371 leaker-based EFs. To convert from leaker EFs to population EFs, GRI-GHGCalc multiplied the 1996 GRI/EPA study Volume 9 leaker EFs by leaks per mile (for mains) and leaks per service (for services) data.

For example, for mains:  $\text{scf CH}_4/\text{hour/leak} * \text{leaks/mile} = \text{scf CH}_4/\text{hour/mile}$

### 2.2. Leak Grading in California

The leak grading criteria used by the state of California[11] follows the GPTC guidelines[12] closely.

Grade 1 leaks require immediate repair or continuous action until the conditions are no longer hazardous. Grade 2 leaks should be repaired within one year, but no later than 15 months from the date the leak was reported. Grade 3 leaks should be reevaluated every 15 months from the date reported until the leak is regraded or no longer results in a reading. SoCalGas terms the

"Grade" of the leak as "Code" rather than "Grade". Leak grading in other states is similar, though not identical, and some states have additional subcategories.

Another difference is that some states say just to monitor grade 3 leaks, whereas California calls on utilities to prioritize fixing high flow Grade (Code) 2s and 3s as soon as possible.

## 2.3. Summary of Studies Referenced or Developed as Part of this Report

SoCalGas collected field leak concentration measurements and leak flow rate data based on a geographically diverse random sample of the entire SoCalGas distribution system, which established a technically sound foundational leak flow rate dataset. A summary of the past industry emission studies referenced, as well as the SoCalGas studies performed as related to this report are summarized in Table 1 below. Further details about the past industry emission studies and SoCalGas studies are shown below Figure 1.

**Table 1: Summary of Studies Referenced and Performed as Part of this Report.**

Study	Header/Legend Abbreviation <sup>(1)</sup>	Year(s) Performed	Report Year	Scale	Number of Samples <sup>(2)</sup>	Used for Report EFs <sup>(3)</sup>	Reference Number
WSU/EDF 2015	Natl_WSU_EDF	2013	2015	National Multi-Utility	212	No	[1]
CARB/GTI 2019	Natl_CARB_GTI	2014-2015	2019	California Multi-Utility	76	No	[4]
OTD/GTI 2013	Natl_OTD_GTI	2011-2012	2013	National Multi-Utility	62	No	[2]
SoCalGas All District Study	AllDisPilot	2019	2019	SoCalGas	78	Yes	This Report
SoCalGas All District Leak Inventory Reduction Program	AllDisLIRP	2019	2019	SoCalGas	10	No	This Report
SoCalGas Decision Tree 157 Pilot	DT157Pilot	2015-2019	2019	SoCalGas	157	Yes	This Report
SoCalGas 3-District Pilot	3DisPilot	2019	2019	SoCalGas	56	Yes	This Report
SoCalGas 3-District Pilot Low Specification	3DisPilotLowSpec	2019	2019	SoCalGas	8	No	This Report

(1) The CARB study was conducted in the state of California, so the inclusion as a "national" study in the analysis is done to allow comparison of two groupings of studies (i.e., past industry studies and SoCalGas studies for this report) without multiple descriptors in table headings and plot legends. The industry studies are all multi-utility, and the SoCalGas studies are only with SoCalGas data. Therefore, in the summary plots and tables of this report, the CARB study is grouped with the two national studies (OTD and WSU below), and the term national studies is retained for this combined set.

(2) - In some cases, a very few (e.g., one or two) site observations were removed from the sample set for comparisons. These were done when upon future dig up of the sites the leak was found to be on a non-pipe item like a valve stem. The individual observations for all studies are in the Appendix section of this report.

(3) All the reports were referenced and statistically analyzed and compared; however, only the three SoCalGas large studies were used to develop the SoCalGas specific EF developed as part of this report. Only the SoCalGas studies included the set of surface concentration measurements as set up and collected using the Decision Tree process.

A basic timeline is provided in Figure 1 below showing the order of activities in developing the Decision Tree process in relation to the various sets of data referenced in the report.

**Figure 1: Timeline of Major Tasks and Milestones.**

Major Tasks and Milestones	2015				2016				2017				2018				2019				2020			
	Q1	Q2	Q3	Q4																				
<b>Development of the DT Process (DT157Pilot dataset)</b>																	◆							
1) Project initiated to evaluate system data and information obtained during leak grading and centering processes																								
2) Go/Nogo Decision - initial study results identified potential approach using CH4 concentration.																								
3) DT process developed and refined. DT157 dataset completed.																								
<b>3 District Pilot (3DisPilot dataset)</b>																	◆				◆			
1) Pilot study kick-off with Operations																								
2) 3 Districts Pilot Study data cut-off for analysis 11/22/19																								
<b>Data Analysis and Technical Report</b>																	◆				◆			
1) Probabilistic Analysis with OTD/Dan Ersoy started																								
2) Draft Report published																								
3) Anticipated publishing data of Final OTD report																								
<b>System Random Sampling (AllDisPilot dataset)</b>																	◆				◆			
1) Random Sampling started																								
2) Random Sampling completed																								
<b>3 District Pilot Lowered Spec (3DisPilotLowSpec dataset)</b>																	◆				◆			
1) Started																								
2) Completed																								
<b>Leak Inventory Reduction Program (AllDisLIRP dataset)</b>																	◆				◆			
1) Started																								
2) Completed																								

## Data Obtained from Past Industry Studies as Comparisons for this Report

**California Air Resources Board.** 2019 Study. (CA\_CARB\_GTI). California, multi-utility (including non-SoCalGas utilities in CA), multi-material, multi-facility study [4]. 76 samples.

**National Operations Technology Development.** 2013 Study. (Natl\_OTD\_GTI). National, multi-utility, multi-material, multi-facility study [2]. 62 samples.

**National Washington State University Study.** 2015 Study. (Natl\_WSU\_EDF). National, multi-utility, multi-material, multi-facility (e.g., service vs. main lines) study [1]. 212 samples.

## SoCalGas Studies for this Report – Data and Study Terminology Definitions

**All District Study (AllDisPilot).** This study is distinct from the other pilot studies noted below. The study sample was stratified by district. A random sample was then drawn with the corresponding number of leak site samples (per district or district grouping) to preserve the correct proportion of the districts in the total SoCalGas population of leaks through all SoCalGas districts. The measurements were performed by GHD Corporation, Air Quality Services Group.

**All District Leak Inventory Reduction Program (AllDisLIRP).** This leak measurement data set consists of leak locations associated with the SoCalGas Leak Inventory Reduction Program (LIRP). The objective was to identify "large" leaks (greater than 10 scfh) within this leak population for

prioritization. The leaks in this group are generally leaks detected many years ago. The locations in this study were system-wide and were measured by company crew or contractors: GHD Corporation & Gas Technology Institute (GTI). A subset of 10 LIRP leaks were provided for the current analysis because those were the only leaks for which concentration measurements and leak flow rate measurements were both obtained. The other LIRP leaks did not have this information and therefore were excluded from the analysis. As will be discussed later, this study was determined to have sample bias based on the sampling criteria because SoCalGas was specifically looking for the *largest* leaks. These 10 samples happened to have surface concentration measurements so they could be reviewed as part of this study.

**Decision Tree 157 Pilot (DT157Pilot).** This leak measurement data set was gathered by the SoCalGas research team to develop the Decision Tree methodology.

**3-District Pilot (3DisPilot).** Sample data from these leak locations were part of an implementation pilot study where the DT process plus subsequent leak flow measurements was deployed in 3 Operating Districts. This data contains the leaks that met the SoCalGas Decision Tree concentration thresholds. As part of the pilot study, all leaks were measured that met at least one of the four possible Decision Tree thresholds. Leaks were measured by company crew or contractors (GHD & GTI). The total screened number of leaks that had surface concentration measurements was 356. Of these 356 leaks, 56 samples had leak flow rate measurements. These three operating districts are considered a good representation of the overall SoCalGas service territory, and they are a geographical subset of the All District Study. The results of an overall sensitivity of leak rate to geographic operating district is presented later in this report.

**3-District Pilot Low Specification (3DisPilotLowSpec).** This includes eight leaks measured within a 3-District Pilot Study area that did not meet the SoCalGas Decision Tree Thresholds for surface concentration percent gas levels that would estimate leak flow rate to be greater than or equal to 10 scfh but were above 1% gas for unpaved surfaces (dirt or grass) or above 5% gas for cracks in paved surfaces. As will be discussed later, this small sample was determined to have sample bias based on the sampling criteria, since leak flow rates were measured at *lower* leak concentration levels for unpaved surfaces and cracks in pavement than the thresholds set by the Decision Tree for the same categories.

## 3. Approach for Field Sampling, Measurement Techniques, and Decision Tree Process

---

### 3.1. SoCalGas Process Development

SoCalGas began to search for a cost-effective means to identify and repair non-hazardous leaks as soon as reasonably possible to minimize the climate change impacts of methane emissions. In 2015, the Research team at SoCalGas began the development of a cost-effective approach of identifying and differentiating leaks on the buried distribution system that have high flow rates, for the purpose of prioritizing repairs.

#### Initial focus on Existing System Data and Leak Centering Process

Initially, the work was based on the hypothesis that existing system leak data obtained at the time buried non-hazardous distribution leaks are detected and graded could be leveraged to identify a sub-set of system leaks that had a greater probability of being high flow-rate leaks. Prior industry studies suggested that the number of system leaks that have a high flow rate is a small percentage of total leaks. The goal was to find an efficient way to identify these leaks, so they can be scheduled for repair. SoCalGas mined existing system data for “leak spread” but found that there was no discernable relationship between the available data and leak flow rate.

The spread of a leak is determined during traditional leak survey identifying the extent of the ground surface area where methane concentration indications are present. This is very useful from a safety evaluation standpoint as it may indicate the potential for hazard to nearby structures. However, as it relates to the flow rate of a leak no discernable relationship was found between the two data sets (see Appendix D).

Next, the concept of relating measurement of emissions from bar holes created during the leak centering process to the actual surface expression leak rate measurements of the leak site was investigated. Over thirty pending leak sites (Code 2s and 3s) were identified for an initial project within the Los Angeles basin area.

The surface expression measurements were completed and documented through detailed field notes indicating the associated concentrations found for each measurement as well as descriptive information of the ground characteristic of each measurement point. Once this was completed, distribution crews were dispatched to drill the bar holes required to center the leak on-site. Measurements of emissions were taken from the centering bar hole created by the crew prior to performing the repairs.

Unfortunately, the centering bar hole concentration data emission rates did not correlate well with the surface expression leak rate measurements of the site. However, some relationships were identified between the surface expression leak rate measurements and the concentration measurement data collected when the field notes indicated ground characteristics were similar.

## Enhanced Leak Survey Practice

As more surface expression leak flow rate data was collected, some relationships of leak flow rate to methane concentration measurements began to emerge. Extensive surface expression method leak flow rate data was then collected along with methane concentration measurement data and corresponding types of ground surfaces and locations. Subsequent groupings of that data based on similar location descriptions yielded promising correlations to the leak flow measurements.

Given that the basis of traditional leak survey practice is to use methane *concentration* measurement data (e.g., % gas) as a means of determining when a leak is present, one should be able to leverage current leak survey practice to provide the additional leak concentration data needed to determine if a leak site has the potential to be a large leak with a higher *flow* rate (e.g., scfh). Another benefit of this approach is that the identification of potential large leaks occurs upon discovery of the leaks, which is in the most ideal timeframe possible.

## 3.2. Development of the Decision Tree Approach

### Types of Surface Conditions

In 2015-2019, SoCalGas visited over three hundred code 2 and 3 leak sites to collect concentration screening value data. Based on the concentration data, 157 locations were selected for surface expression flow rate measurements (DT-157 SoCalGas Pilot).

The surface expression measurements were analyzed against site-specific characteristics (e.g., unpaved, crack in pavement, etc.) of ground-level methane concentration measurements. Next, concentration thresholds were developed for each surface category to identify the leaks with the potential for a high flow leak rate measurement.

A process was then developed where surface expression measurements to calculate leak flow rate are performed whenever one or more of the leak concentration threshold values are met or exceeded, so repairs can be prioritized.

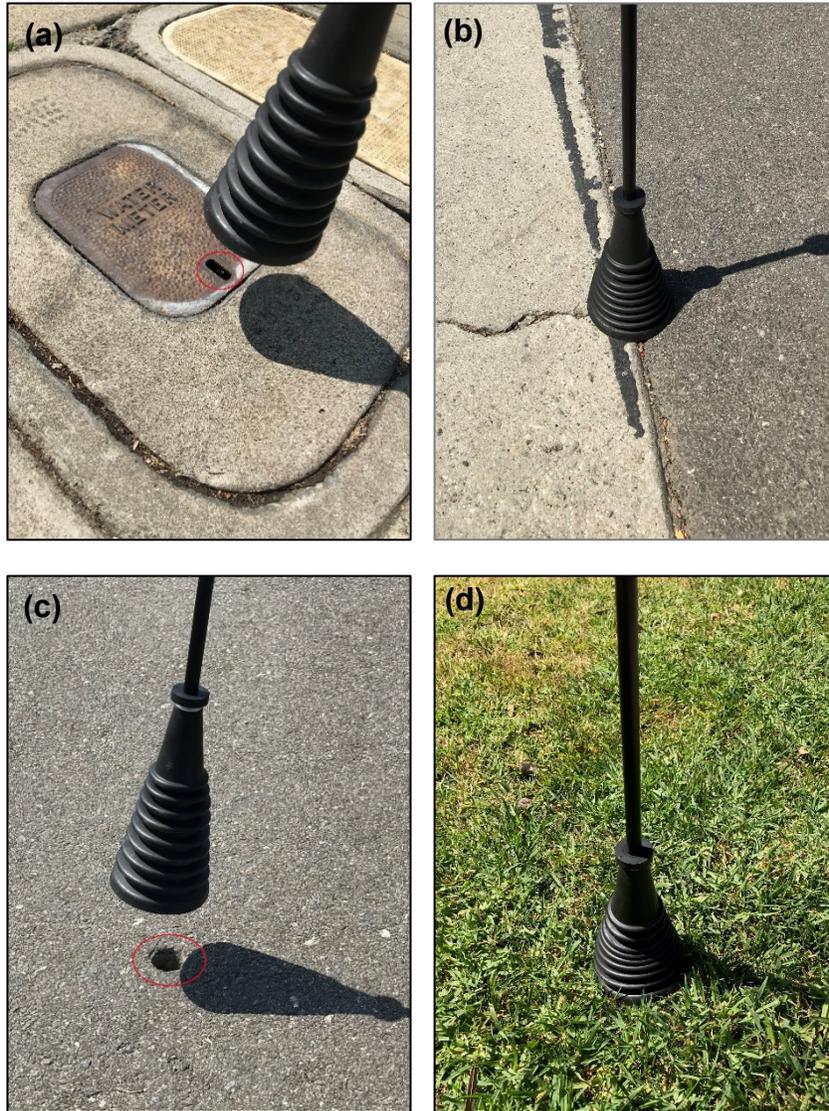
This methodology was termed the “Decision Tree” (DT) approach.

The DT approach collects the *maximum* methane concentration measurements at defined types of surface condition locations. The current process, which was also used to collect data as part of this study, uses four defined types of surface condition locations:

- Crack (or seam) In Pavement - CIP
- Unpaved Surface - US
- Bar Hole (leak survey type) - BH
- Small Sub-Structure (not gas system related) - SSS

The defined types of surface condition locations are shown in Figure 2 below.

**Figure 2: Surface Condition Location Definitions.**



(a) Sub-Structure-SSS, (b) Crack in Pavement-CIP, (c) Barhole-BH, and (d) Unpaved Surface-US

### **Methane Concentration Measurement Process**

The operator utilizes the Heath DP-IR (Detecto Pak Infrared) or GMI Gasurveyor along with the survey probe attachment to survey the leak site by placing the conical end of the survey probe (as shown in Figure 2 above) directly onto the ground surface collecting drawn samples of methane in air.

This device contains a pump that draws the air samples from the cone-shaped probe at the ground surface to an analyzing chamber where infrared lasers are used to quantify the concentration of methane in the air and provide a methane concentration reading to the operator.

The spread of the leak is determined by probing the ground sub-surface and identifying the extent at which elevated methane concentrations are present.

Once the spread is determined the operator then identifies and records the highest sustained reading within the spread of the leak for each of the four surface conditions.

The required placement of the cone shaped probe for each of the 4 surface conditions is as follows:

- Crack (or seam) In Pavement - CIP: The probe is placed directly on top of the crack or seam and is in contact with the paved surface.
- Unpaved Surface – US: The probe is placed directly on top of and in contact with the unpaved surface whether it be soil, grass, or rocks.
- Bar Hole (leak survey type) – BH: The probe is placed directly on top of the bar hole with the cone shaped probe in contact with the ground surface and encompassing the bar hole.
- Small Sub-Structure (not gas system related) – SSS: The probe is placed directly on inside of or on top of the access hole(s) of a substructure (prior to venting or lifting the lid) with the cone shaped probe in contact with the ground surface and encompassing the access hole(s).

Training was provided to the operators in the pilot study to familiarize them with these additional requirements.

## **Decision Tree Concentration Threshold Values for Leak Flow Rate Measurement**

The applicable *concentration* measurements are then compared to the threshold values, and any one (or more) of the up to four concentration measurements that meets or exceeds the threshold concentration value will then result (trigger) in that leak being classified as a potentially large, non-hazardous leak with a predicted leak flow rate of 10 scfh or higher.

The threshold values for concentration measurements by surface condition type are:

- 20% gas: Crack (or seam) In Pavement - CIP
- 5% gas: Unpaved Surface - US
- 80% gas: Bar Hole (leak survey type) - BH
- 60% gas: Small Sub-Structure (not gas system related) - SSS

These leak sites that trigger at least one of the defined types of surface conditions are then scheduled for leak flow rate testing using the process described below.

## **Leak Flow Rate Measurement Process**

Leakage *flow* rates were measured from selected underground distribution pipeline leaks, chosen based on the Decision Tree methodology. The leakage flow rates were measured using the well-established and published Surface Expression methodology (detailed in Appendix A). Leak flow rates are reported in scfh of methane (CH<sub>4</sub>).

These leak rate measurements provide an approximation of ‘in-air’ methane emission rates without the need to excavate the leak source. The study raw field data, including both concentration measurements and leakage flow rates, are in Appendix E.

### 3.3. Selection of Single vs. Facility/Material Specific EFs

SoCalGas also evaluated leakage spread data and “material” (e.g., plastic vs. steel, etc.) and “facility” (e.g., service vs. main) data to determine whether any of these system variables would help in predicting large leaks.

The discovery of the common mis-association of predicted material and facility data resulted from this effort and is demonstrated in the SoCalGas results of a collaborative study with CARB [4]. Based on this limitation, a single emission factor approach was selected.

### 3.4. Precision and Sample Size Analysis - Minimum Sample Size

A precision and sample size analysis (similar to a power and sample size analysis for hypothesis analysis) for desired confidence intervals was conducted ahead of the sampling[13, 14].

Precision and sample-size (PrSS) analysis is a key component in designing a statistical study that uses confidence intervals (CIs) for inference. It investigates the optimal allocation of study resources to increase the likelihood of the successful achievement of a study objective.

There is a strong correspondence between CIs and hypothesis tests. A  $100(1-\alpha)\%$  CI can be obtained by inverting the acceptance region of the corresponding level  $\alpha$  test. In other words, a  $100(1-\alpha)\%$  CI provides the entire range of hypothetical values for a parameter of interest that cannot be rejected by the test at a significance level of  $\alpha$ .

Despite the strong correspondence between PrSS used for CI analysis and Power and Sample Size (PSS) used for hypothesis tests, they will not necessarily lead to the same requirements for the sample size. A hypothesis test compares the parameter of interest with a single value, whereas a CI provides a range of plausible values. Thus, for the same significance level, the sample-size requirements for the CI will generally be larger than for the hypothesis test.

#### Sample Size for a One Mean Confidence Interval

Although SoCalGas used the bootstrap method to infer these population mean leak rates, the classical method of Precision and Sample Size is a useful tool while setting up a sampling plan. The selection of  $1-\alpha$ , probability of achieving the CI, and desired precision (i.e., the width of CI) are subjectively set; therefore, multiple values for these parameters are often selected for comparison. The same can be true of the estimated standard deviation of the measure of interest.

Using the typical/expected values from the national studies and the SoCalGas DT-157 studies, a set of inputs to the analysis was established:

- Confidence Level (level): 95%

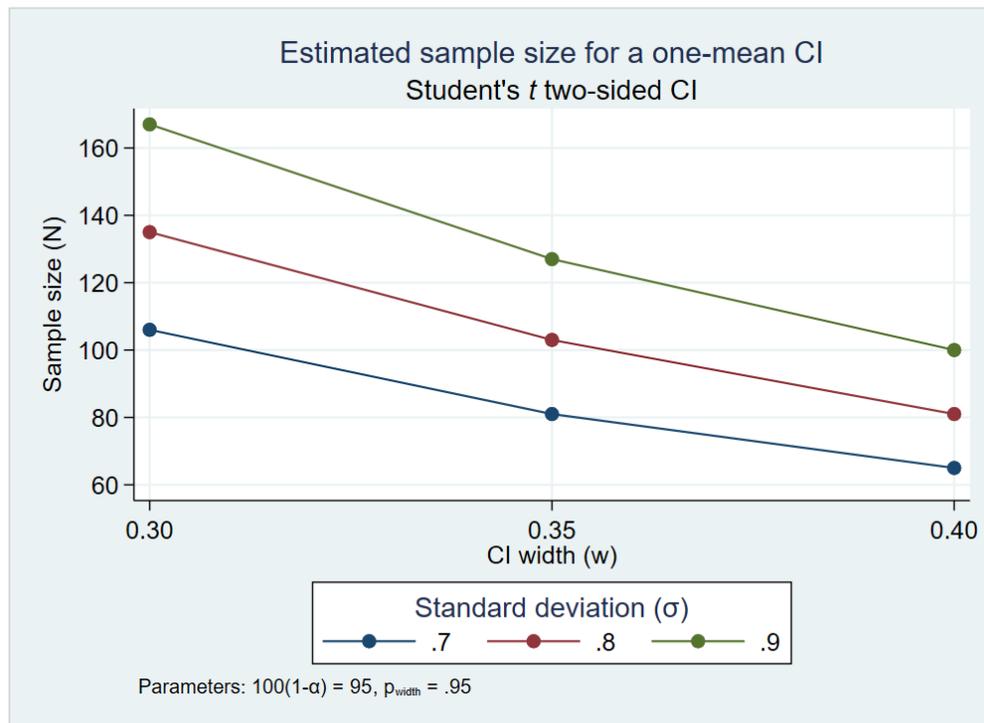
- Estimated standard deviation (SD) of log(10) of the Leak Rate in scfh: 0.7 - 0.9
- Desired precision (width) of log(10) of leak rate in scfh: 0.3 - 0.4
- Desired probability of achieving the CI width (Pr\_width): 95%

With these input values, a sample size (N) calculation for a one mean CI of the log(10) of the leak rate was conducted and the results are presented in both Table 2 and Figure 3 below.

**Table 2: Sample Size for a One Mean Confidence Interval of Log(10) Leak Rate.**

Level	N	Pr_width	width	sd
95	106	0.95	0.3	0.7
95	135	0.95	0.3	0.8
95	167	0.95	0.3	0.9
95	81	0.95	0.35	0.7
95	103	0.95	0.35	0.8
95	127	0.95	0.35	0.9
95	65	0.95	0.4	0.7
95	81	0.95	0.4	0.8
95	100	0.95	0.4	0.9

**Figure 3: Sample Size for a One Mean Confidence Interval of Log(10) Leak Rate.**



A centered value of 100 samples is shown in Table 2, which is a good target sample size to collect. This study resulted in 291 samples, which meets the one-mean sample size requirement for a standard deviation of 0.9 and CI width of 0.3 at 95% confidence.

## Sample Size for a Bayesian Proportional Analysis

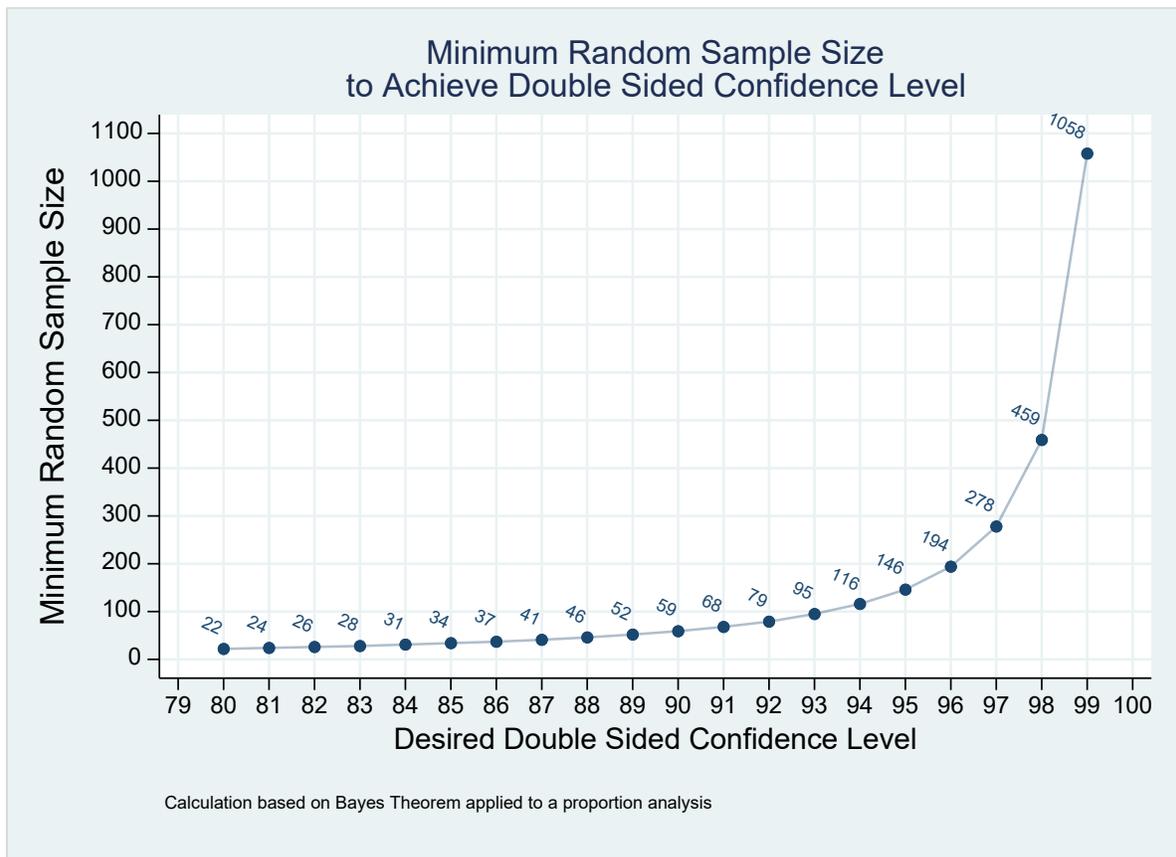
A similar analysis for minimum samples sizes was conducted but related to the Bayesian proportional analysis, which was done to create the Decision Tree performance error table.

The term proportional, as related to Bayesian analysis, refers to an analysis that is measuring the proportion (e.g., percentage) of a sample that falls in a certain category vs. calculating a variable of interest and its expected value and any other parameters. In this case, the Bayesian analysis is coherent and provides the expected value as well as the credible intervals at a 95% confidence level for the errors in the DT predictions.

Figure 4 below shows the minimum sample size required for a given two-sided (two tailed) confidence level. This is the sample size for any category of leak flow rate that will bring the Bayesian analysis Upper Prediction Limit (UPL) below the specified single sided confidence limits for the analysis when there is no occurrence in that category of a leaking sample.

The calculation for minimum sample size uses a conservative uniform (aka ignorant) prior for the Bayesian analysis to establish this significant sample size[15]. The sample size was established to ensure a better than 95% single sided credible band (aka, upper or lower prediction limit) was achievable for the proportion analysis.

**Figure 4: Sample Size for a Bayesian Proportional Analysis by Confidence Level.**



Ultimately, the error analysis had 117 (for Type I error calculation) and 174 (for Type II error calculation) respectively, comfortably above the desired size. This established a conservative number of samples, since a uniform distribution is used for the proportion prior. This will be described in a later section of this report.

### 3.5. Physical Sampling Plan - Where to Sample

It is essential that a simple random or stratified random (for efficient, smaller sample sizes) sample be used vs. any other, potentially biased sample. A random sample is designed so each leak has the same likelihood of being chosen as part of the sample from the population of leaks. It will improve the accuracy of the results, helping ensure the statistics inferred from the sample are representative of the population. When such a proper random sample is taken then the analysis done from the sample such as bootstrap, Monte Carlo, regression, Bayesian proportions and the inferred results from the same are statistically sound.

For the SoCalGas *All District Study*, the non-hazardous leak population was stratified by operating district. A random sample was then drawn with the corresponding number of leak site samples (per district or district grouping) to preserve the correct proportion of leaks in each. This study contained 78 samples, which as discussed earlier and can be seen in Figure 4 is greater than the 65 samples required to achieve the desired 95% double-sided credible band with a CI width of 0.4 and standard deviation of 0.7.

As was done in this study, during the development of the Decision Tree process and its ongoing application for leak measurements, one should also take a random sample from the leaks that the Decision Tree estimates as *not* large leaks (i.e., less than 10 scfh) to confirm for true/false negatives. This is important to: (a) continually scan the population for situations that may have changed over time, and (b) continue to confirm and refine the Decision Tree performance metrics.

## 4. Methodology Overview of Data Collection and Statistical / Probabilistic Analysis

---

### 4.1. Descriptive Statistics of Study Samples

Standard descriptive statistical analysis was conducted including calculation of sample means, medians, percentiles, interquartile ranges, and other statistics. Various analysis and plotting techniques were used to confirm sampling bias and draw high-level conclusions on the different individual and grouped sample leak rate distribution center tendencies, uncertainties, and shape. An evaluation of SoCalGas data alongside recent Industry studies served as a baseline comparison.

### 4.2. Data Transformation, Regression, and MCMC Models

Data transforms were used to ensure that any regression utilized had a sound basis. Linear Regression (LR), and Analysis of Variance (ANOVA) were used to quality check sample set data, spot outliers, confirm assumptions, assess regression, and assess probabilistic residuals and diagnostic measures.

These steps led to the selection of a non-biased and robust SoCalGas sample set for further analysis. The highlights of these analysis methods are presented in the report body with additional, supporting regression details in Appendix B.

### 4.3. Decision Tree Predictive Performance

A purely probabilistic Bayesian analysis[15-17] was used to measure the Decision Tree performance that grouped detected leaks into two categories: “Large” leak rate (greater than or equal to 10 scfh) and “Not Large” leak rate (less than 10 scfh) predicted leak *rate category* based on site *concentration* measurements. A full false/true negative/positive error table (Type I and II error) was developed with Bayesian derived lower and upper prediction (credible) limits.

### 4.4. Population Mean Leak Rate Inferential Analysis

Resampling with replacement (Bootstrap) analysis[18-20] of field leak rate data and Monte Carlo sampling of a fitted data distribution of leak rates were both used to infer the population mean leak rates with upper and lower confidence limits from the sample data.

### 4.5. Emission Factor Determination

The SoCalGas emission factors were derived using the combination of the appropriate bootstrap population leak rate means and the Bayesian Decision Tree error table percentiles.

## 4.6. Method of Emission Factor Application

Several scenarios were developed to show how a utility could apply the emission factors based on operational considerations and if leak rate measurement and/or concentration measurements were taken.

## 4.7. Quality Assurance

### Statistical Checks

The field data sample sets were reviewed from a data quality perspective through statistical quality checks, including:

- To establish if there were unexplainable outliers or extreme values, a set of statistical diagnostics was conducted, including: DFBETAs, Lowess, Leverage, and others.
- Normality of the response variable (leak rate) and the associated predictive residuals were checked, leading to applying a log transformation of the leak rate data which achieved a normal distribution of the transformed data.
- A combination of ANOVA, Regression, and non-parametric statistical tests (e.g., Kolmogorov-Smirnov) were conducted to test for normality of variable distributions, comparison of sample means, and other statistical parameters. These tests are noted in the report and/or documented in the appendices when appropriate.
- Since there is a danger in over relying on simple to use and quantitative statistical tests, this research also examined the normal quantile plot to determine normality rather than blindly relying on a few test statistics[21, 22]. For example, the Kolmogorov-Smirnov test[23] generally is not very powerful against differences in the tails of distributions. For these reasons, the Kolmogorov-Smirnov is not a particularly powerful test in testing for normality[24]. Hence, the quantile normal plots of the data were also carefully analyzed.

### Probabilistic Regression Check

- As a secondary check on the traditional linear regression and analysis of variance (ANOVA), a Bayesian Monte Carlo Markov Chain (MCMC)[25] Metropolis-Hastings Sampling (MHS)[26, 27] random walk and Gibbs Sampling (GS)[28] non-frequentist and non-parametric analysis was also completed.

### Significant Figures

- The values of leak flow rate measurements are reported to an extended number of digits in tables to allow comparison and prevent cumulative rounding errors when conducting calculations involving multiple variables and bootstrap analysis. These values are typically listed with three *digits* after the decimal place.
- The limiting significance values come from the precision of the flow rate measurement equipment (scfh), and the calculated emission factors (EF) are limited to two digits based on this precision.

- When the average leak flow rate data is combined with the Bayesian expected proportions for the Decision Tree assignments to calculate Emission Factors (EF), the total number of significant figures reported is three, which results in two digits after the decimal place.
- The output of regression analysis is automated and reported to many decimal places in the standard regression output table. These values were retained, and it should be noted that the precision of the analysis is not represented by these formats.

### Standard Conditions

- The standard conditions related to the leak concentration and leak flow rate measurements in the report are to 1 atm of pressure and 60F for temperature.
- Additional details on temperature compensation of equipment is provided in the appendix.
- Leak flow rate measurements and the emission factors derived from the report are reported in scfh methane (CH<sub>4</sub>).

### Distribution Fitting and Monte Carlo Sampling

- A series of distribution fits was conducted on the log(10) of the leak rate for the SoCalGas sample set (291 samples). The results are described in Appendix C and include a few of the selected distributions that were fit and their goodness of fit measures.
- For illustrative purposes, the log-normal distribution fit was selected. The log-normal fit was sampled using a Monte Carlo method where the sample size of the distribution was set to the same sample size used to fit the distribution. This ensures that the average leak flow rates estimated from the Monte Carlo analysis contain the uncertainty associated with the limited sample size from which they were derived.
- As noted, the distribution fit was checked with multiple goodness of fit parameters, and the error associated with the fit done by the associated software was orders of magnitudes smaller than the uncertainty associated with the random sampling of the distribution (as noted below), so this fit error was not considered.
- The log-normal distribution fit was *not* used to calculate the emission factors, since a full bootstrap resampling with replacement analysis was done on the actual field leak flow rate sample measurements. However, the two were compared in the report to illustrate that if an operator does not have an adequate sample to run a bootstrap analysis of the average leak flow rate, then sampling from a fitted distribution could provide a "stop-gap" alternative until a large enough sample size from the field is obtained.

### Leak Concentration and Flow Rate Measurement Error/Uncertainties

The concentration measurement and flow rate measurement uncertainty for the techniques used in this study have been laid out extensively in the referenced reports. In this study, the leak flow rate measurements are considered the baseline standard (i.e., an accurate indication of the true leak flow rate), and the concentration measurements are used to trigger the Decision Tree threshold points to determine if a leak is assigned to be a predicted *Large* or *Not Large* leak.

Therefore, the Bayesian analysis which calculates the true/false negatives and positives of the DT assignments inherently includes all measurement errors and uncertainty which is folded into those proportions. Further, the uncertainty from the population bootstrap resampling is many times larger than the measurement error, as is the uncertainty generated from the Bayesian proportional analysis and the associated upper and lower prediction limits.

### **Carrying Uncertainty Through to the Emission Factor Calculations**

Additional steps are necessary to properly carry through the uncertainty related to the average (i.e., expected or baseline) emission factor and provide confidence limits at a selected confidence level for the EFs.

To do this, one would run Monte Carlo analysis by drawing from the bootstrap average leak rate population distributions of the appropriate data set and category of the leak rate (large and not large) and then weight those by the Bayesian proportions for those categories. This would be done thousands of times, picking the average leak flow rates and the associated Bayesian proportions from those distributions and then calculating the associated emission factors.

This would provide a full distribution of the emission factors for each category; one could then select the confidence level of choice (e.g., 95%) to generate the confidence interval around the average emission factors.

However, one would still use the expected (average) value of the emission factors in practice, but the confidence bands would help establish the level of uncertainty in those values.

SoCalGas plans to continue to implement the DT and leak flow measurement process system-wide and collect additional samples from ongoing leak surveys. This greatly increased data set will eventually be used to map the full uncertainty through the entire process to allow a set of confidence bands to be calculated for the emission factor's expected values, i.e. the base case. As of now, it should be noted that this information is not available to the emissions estimates currently being reported by the industry.

## 5. Analysis and Results

---

### 5.1. Descriptive Statistics

#### Comparison of SoCalGas Study to Industry Studies

The data obtained from the eight studies (described above) were quality checked and arranged into a single flat table.

Each of the studies used a Hi-Flow sampler and the dynamic flux chamber method. The main differences are likely the shape/material of the enclosure and the CGI used at the outlet of Hi-Flow. The minimum quantifiable leak rate, assuming a CGI accuracy of 5 ppm is approximately 0.002 scfh.

- WSU/EDF 2015: Hi-Flow + enclosure + CGI
- CARB/GTI 2019: Hi-Flow + enclosure + CGI
- GTI/OTD 2013: Hi-Flow + enclosure + CGI
- Five SoCalGas Studies: Hi-Flow + enclosure + CGI

As noted in a previous section, three Industry leak rate studies were used as baseline studies to which the SoCalGas pilot study were compared and contrasted.

Additionally, there were five SoCalGas studies. Three of these studies are considered "core" studies and are the focus of the emission factor calculation. These include the 3-District Pilot, the All District Study, and the Decision-Tree (DT) 157 study.

The remaining two SoCalGas studies (3DisPilotLowSpec and AllDisLIRP) contained very limited sample sizes of 8 and 10 samples respectively. In addition to being limited, these sample sets exhibited known sampling bias as will be demonstrated later in the report. That said, these data sets are still included in the analysis for comparison as well as to demonstrate how a biased sample could impact the statistical analysis and therefore bias the resulting emission factors.

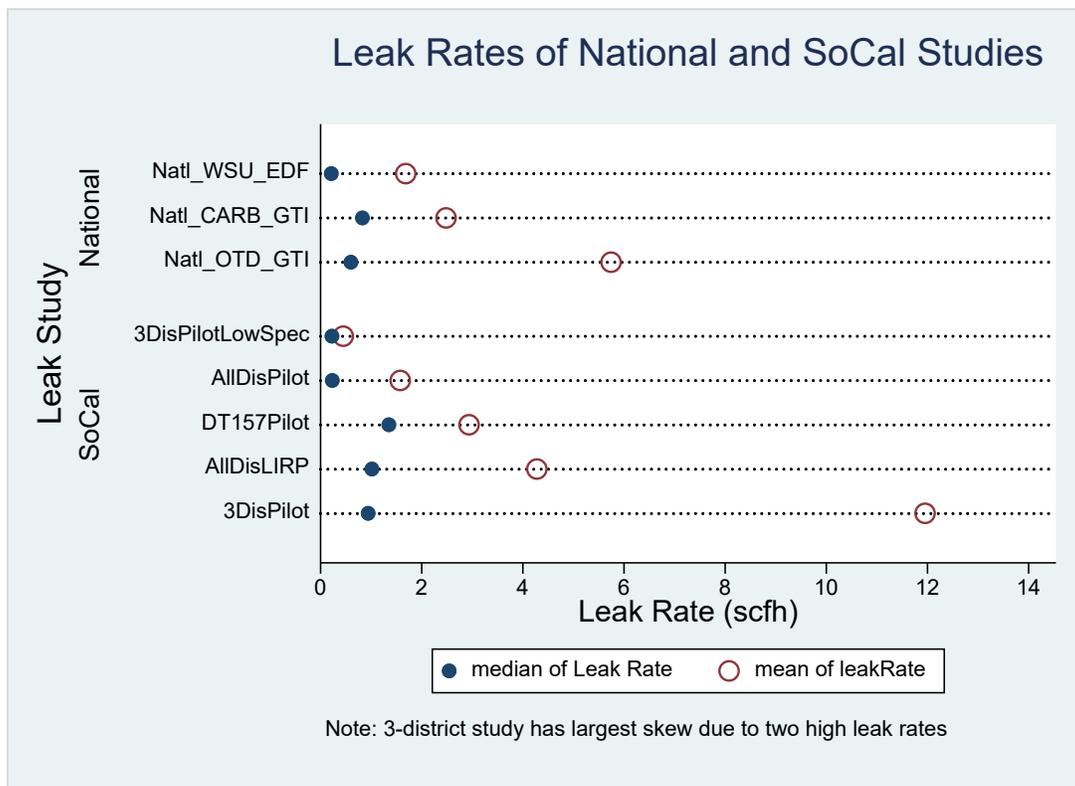
Sample counts as well as *methane* leak rate mean, minimum, and maximum for each of the eight studies are compiled in Table 3 below.

**Table 3: Leak Rate Mean, Minimum, and Maximum by National and SoCalGas Study.**

National				
Study	N(count)	mean(scfh)	min(scfh)	max(scfh)
Natl_CARB_GTI	76	2.481	0.007	20.400
Natl_OTD_GTI	62	5.743	0.044	95.400
Natl_WSU_EDF	212	1.683	0.003	109.472
<b>Total</b>	<b>350</b>	<b>2.576</b>	<b>0.003</b>	<b>109.472</b>
SoCalGas				
Study	N(count)	mean(scfh)	min(scfh)	max(scfh)
3DisPilot	56	11.952	0.020	373.000
3DisPilotLowSpec	8	0.448	0.060	1.640
AllDisLIRP	10	4.276	0.192	30.702
AllDisPilot	78	1.575	0.003	27.045
DT157Pilot	157	2.935	0.003	43.776
<b>Total</b>	<b>309</b>	<b>4.205</b>	<b>0.003</b>	<b>373.000</b>
<b>Grand Total</b>	<b>659</b>	<b>3.340</b>	<b>0.003</b>	<b>373.000</b>

The mean and median of the eight studies are plotted in Figure 5.

**Figure 5: Leak Rate Median and Mean Plot by National and SoCalGas Study.**



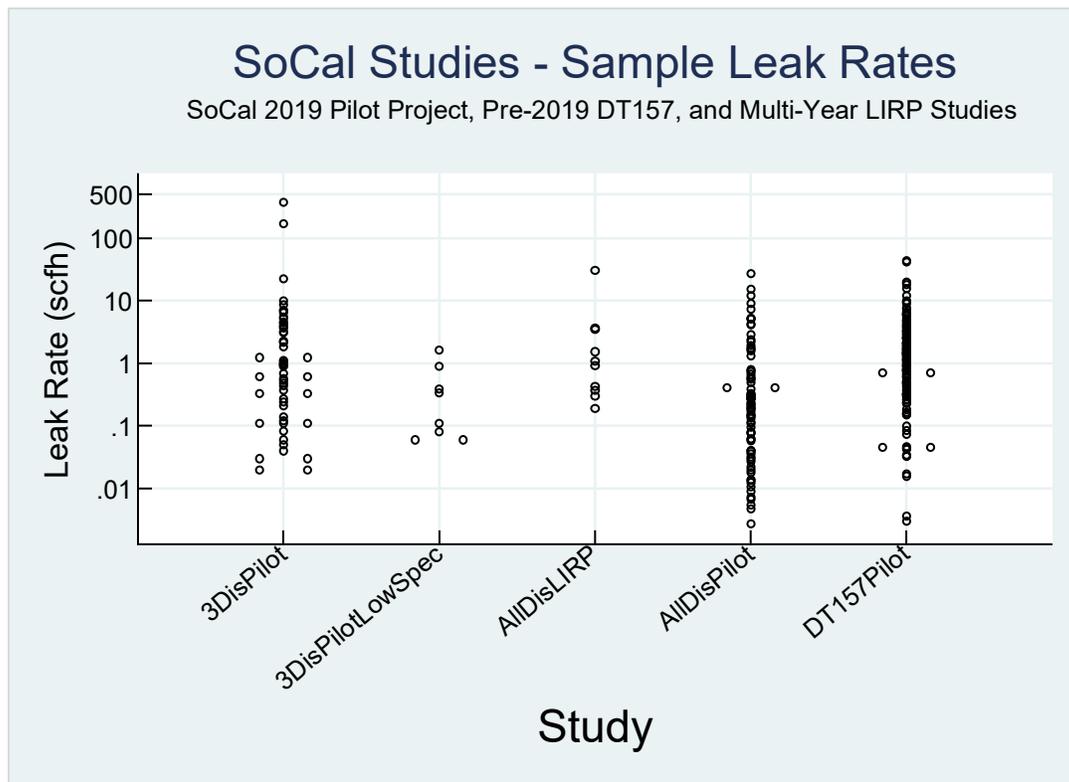
### Data analysis using Dot Plots

The individual leak samples are plotted in a vertical dot plot and delineated by SoCalGas study (Figure 6). This type of plot is similar to a two-variable scatter plot, but, for the purposes of this analysis, the dot plots are used to visualize single variable trends of data. Most data for the three studies fall between 0.1 and 10 scfh except for the two smaller studies where 3-District Pilot Low Specification does not have leak rates *above* 1 scfh except for one point. Additionally, the All District LIRP data does not have leak rates *below* 0.1 scfh. The All District Study also has significantly more data on the lower end of the leak scale below 0.1 and 0.01 scfh.

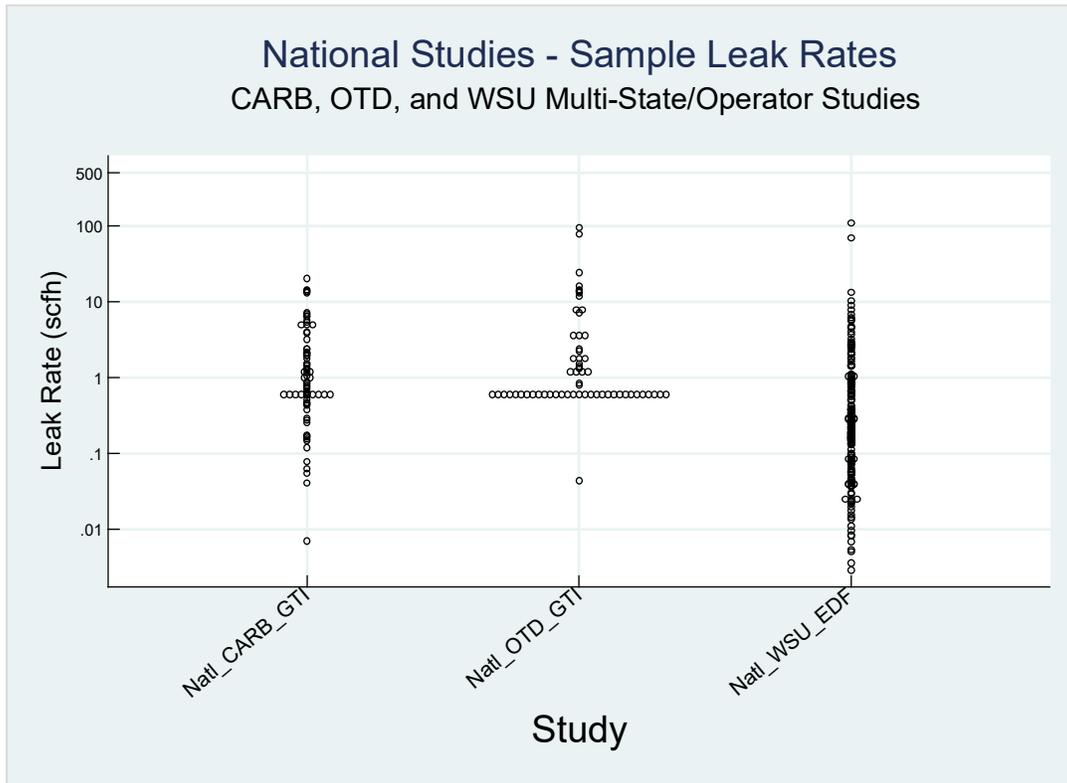
In general, the various studies have a relatively large spread, but the 3-District Pilot Low Spec and the All District LIRP have a tighter spread, which may be a result of a smaller sample set as shown in Table 3 or another factor such as sampling bias. A further analysis of these two studies is presented later in this report.

A vertical dot plot of the three national studies is shown in Figure 7 below which reveals that the WSU study has significantly more data with leak rates less than 0.1 and 0.01 scfh. This may reflect the lower detectable leak rate in the WSU study (0.003 scfh) compared to the other national studies.

Figure 6: Leak Rate Plot by Sample for Five SoCalGas Studies.



**Figure 7: Leak Rate Plot by Sample for Three National Studies.**



The leak rate means and 95% confidence interval limits (5% and 95%) (grouped by national and SoCalGas studies) are shown in Table 4 below. The confidence intervals show similar lower levels, but the combined data set for all five SoCalGas Pilots has a higher upper level as would be expected by the higher maximum values in the data set and their distribution. Focus will be placed on the bootstrap confidence intervals later in this report as they are more robust against issues with non-normally distributed variables.

**Table 4: Leak Rate Mean and 95% C.I. National vs. SoCalGas Studies.**

Study	Obs	Mean	[95% Conf. Interval]	
National	350	2.576	1.534	3.617
SoCal	309	4.205	1.551	6.860
<b>Total</b>	<b>659</b>	<b>3.340</b>	<b>1.980</b>	<b>4.699</b>

The sample count and sample percentiles 5%, 50% (median), and 95% are compiled in Table 5 below.

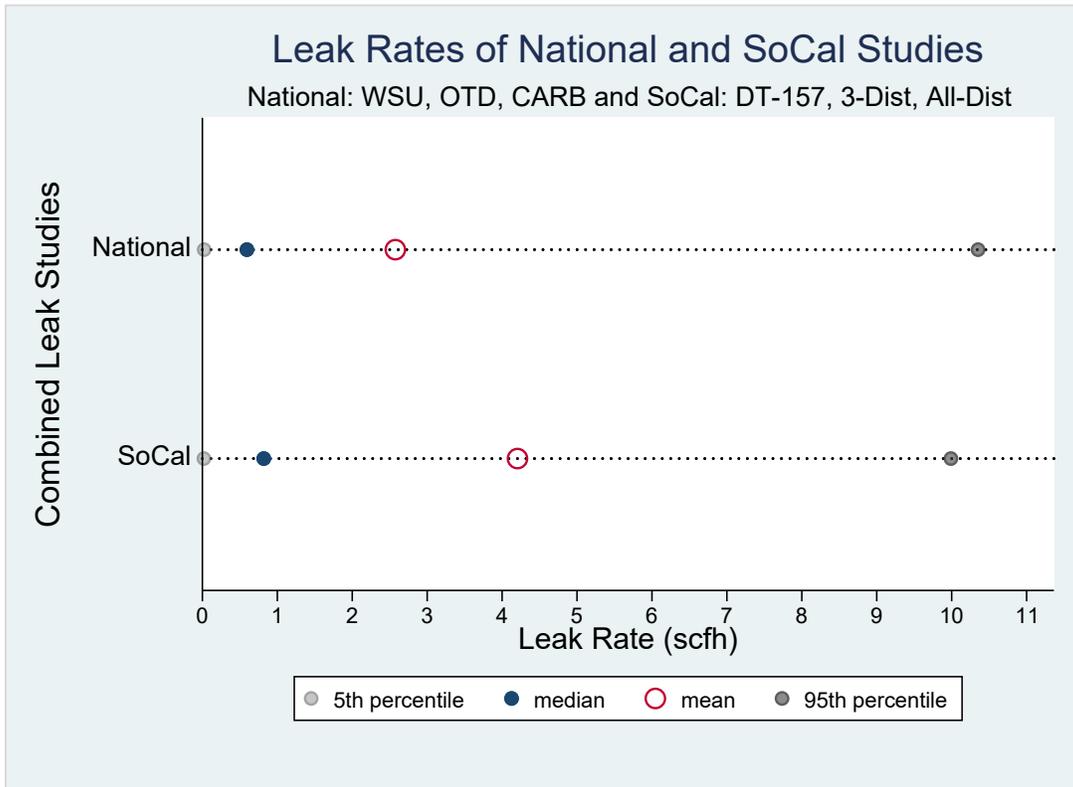
**Table 5: Leak Rate Median, 5%, and 95% Percentiles by National and SoCalGas Study.**

<b>National</b>				
<b>Study</b>	<b>N(count)</b>	<b>p5(scfh)</b>	<b>med(scfh)</b>	<b>p95(scfh)</b>
Natl_CARB_GTI	76	0.063	0.827	13.800
Natl_OTD_GTI	62	0.600	0.600	16.200
Natl_WSU_EDF	212	0.014	0.211	4.753
<b>Total</b>	<b>350</b>	<b>0.024</b>	<b>0.594</b>	<b>10.352</b>
<b>SoCalGas</b>				
<b>Study</b>	<b>N(count)</b>	<b>p5(scfh)</b>	<b>med(scfh)</b>	<b>p95(scfh)</b>
3DisPilot	56	0.030	0.940	22.290
3DisPilotLowSpec	8	0.060	0.225	1.640
AllDisLIRP	10	0.192	1.015	30.702
AllDisPilot	78	0.007	0.231	9.192
DT157Pilot	157	0.042	1.350	9.990
<b>Total</b>	<b>309</b>	<b>0.020</b>	<b>0.819</b>	<b>9.990</b>
<b>Grand Total</b>	<b>659</b>	<b>0.020</b>	<b>0.600</b>	<b>10.000</b>

A simple but compelling plot of the key points of the two tables above is shown in Figure 8 below. This plot shows the combined datasets from national and combined SoCalGas studies median, mean, 5<sup>th</sup> percentile, and 95<sup>th</sup> percentile of the samples sets.

There is especially strong similarity to the percentile (5%, 50%, and 95%) statistical measures between the two large combined data sets of 0.024 vs. 0.020, 0.594 vs. 0.819, and 10.352 vs. 9.990 respectively.

**Figure 8: Leak Rate Lower/Upper Percentiles by Combined National and SoCalGas Studies.**



A similar plot to Figure 8 above, but with all eight studies will be presented in an upcoming section on sample bias and will be discussed in association with identification of studies with sample related bias.

### Data analysis using Box Plots

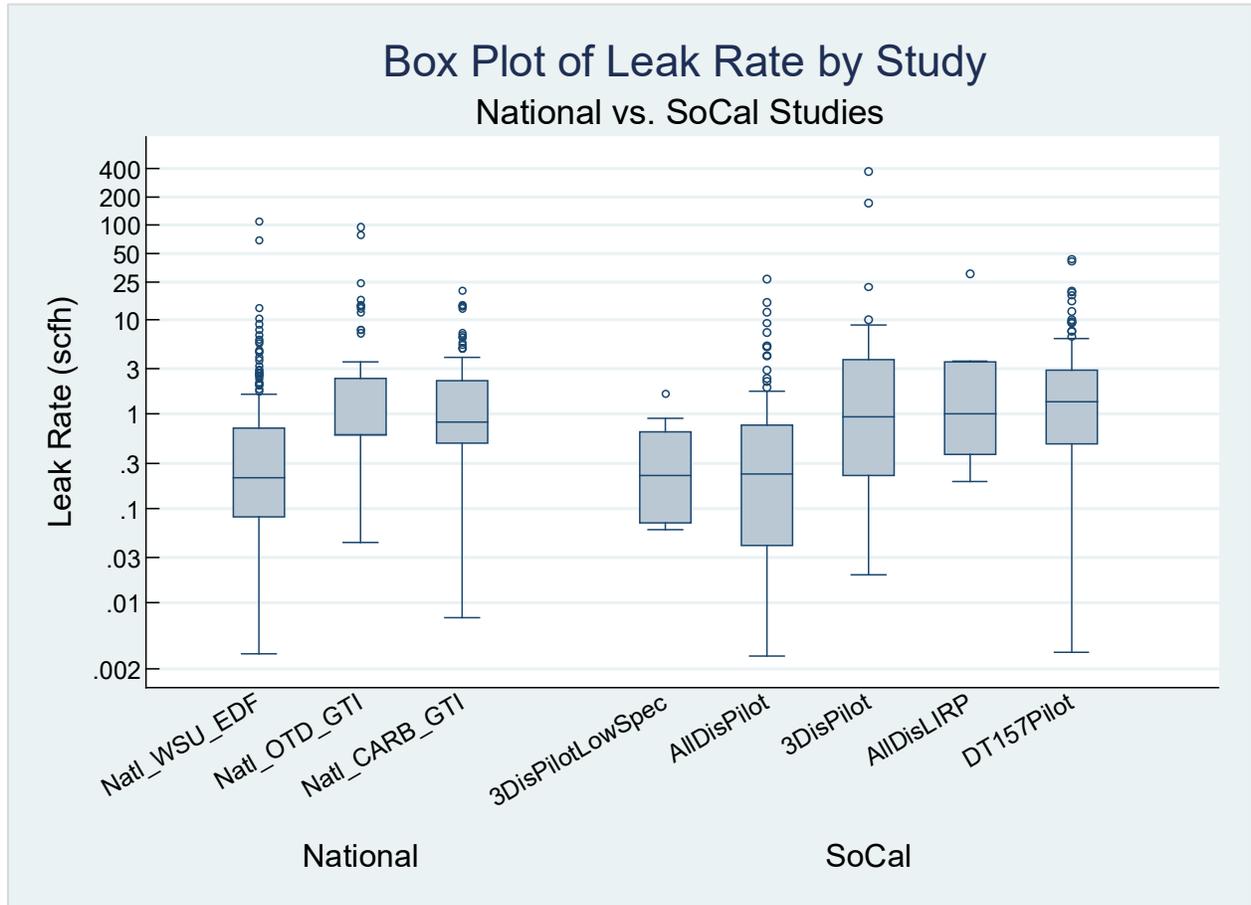
A box plot of the same leak data from the eight studies is shown in Figure 9. The box plot contains 50% of the data within the box, 25% above and 25% below the box. The median is drawn as a horizontal line within the box. Half the data is above and below the median. The box tails (whiskers) extend up to lesser of either the most extreme value or 1.5x the box height (also known as the inner quartile region, IQR). Values outside the tails are known as outliers or extreme values in general - but this should not infer that there is something wrong with the data, rather a close look should be taken of these points. As both the vertical dot and box plots are using logarithmic scales for the y-axis (leak rate), one can see that the distribution of the data appears to be log-normally distributed. This will be further analyzed in the Transformations section of this report.

The national OTD and CARB studies have similar distributions of data with nearly overlapping IQRs. However, the WSU study has nearly the entire IQR below the CARB and OTD studies. For the SoCalGas studies, the same trends (as explained in the vertical dot plot) are apparent with the 3-District Pilot Low Specification and All District Study. The studies have noticeably lower IQRs,

with the 3-District Pilot having its highest value within the upper whiskers (or IQRs) of the other studies. This indicates low bias.

Additionally, the All District LIRP study shows its lowest value within the IQR's of the other studies, indicating high bias. This will be more quantifiably reviewed using frequency plots in the next section.

**Figure 9: Leak Rate Box Plots by National and SoCalGas Studies.**



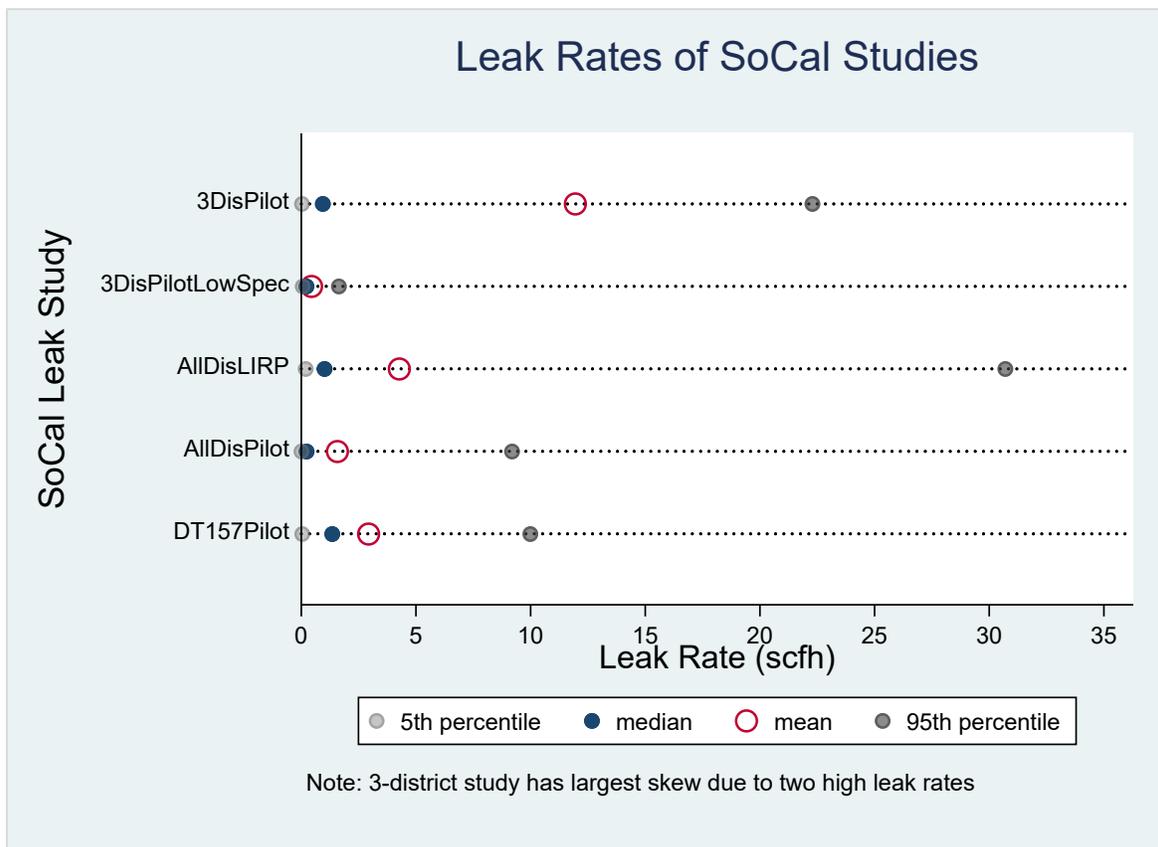
## 5.2. Removal of Studies with Sample Bias

Based on the analysis in the section above and more critically the knowledge of "confirmation bias" for the sampling of the 3-District Pilot Low Specification and the All District LIRP, these two studies will be removed from the three core studies prior to further analysis and incorporation into the emission factor calculations.

The five studies are plotted on a horizontal dot plot in Figure 10 below to allow comparison of study means, medians, and lower 5<sup>th</sup> and upper 95<sup>th</sup> percentiles. The 3-District Pilot Low Specification study has the lowest mean, median, and 95<sup>th</sup> percentile of the five SoCalGas studies. This study's sample selection used a lower Decision Tree concentration criterion (see Approach section for the criteria values) to check on the performance of the criteria and to look for additional false negatives by triggering a leak rate measurement at a lower threshold of leak concentration values. One would expect this to bias the sample to lower values than a purely random sample, which is the case.

The same is true of the All District LIRP study where this sample has the second highest mean and the highest 95<sup>th</sup> percentile. This study's sample selection was designed to pick higher rate leaks for the study. One would expect this to bias the sample to higher values than a random sample, which is the case. A look at the cumulative fraction plots will reinforce these observations.

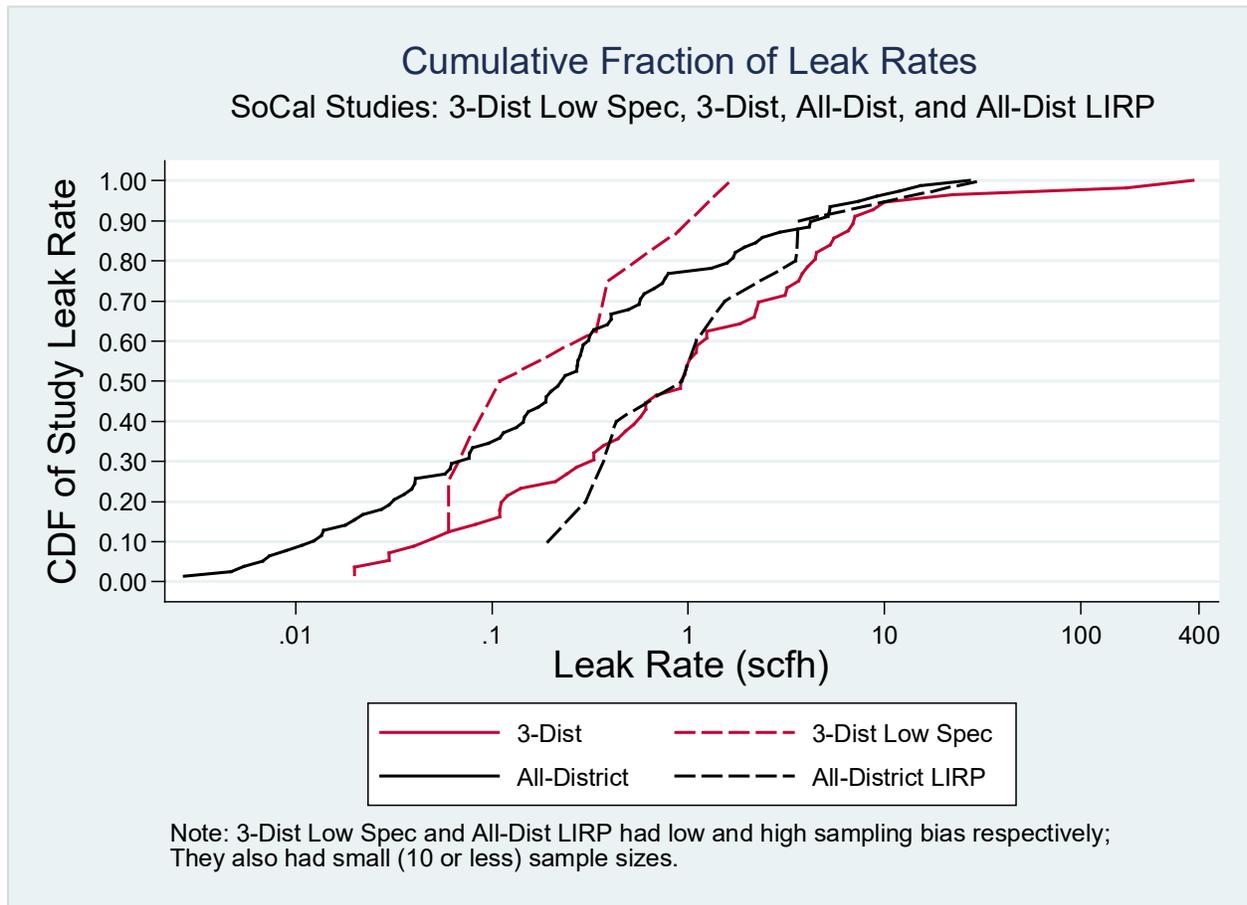
Figure 10: Leak Rate Lower and Upper Percentiles by SoCalGas Study.



The effect of sampling bias on the two smaller studies is shown below in Figure 11 where the biased studies (dashed lines) are shifted to the left (3-District) or right (All District) of their associated larger and randomly sampled studies.

The remainder of this report (with exception of one general ANOVA analysis) will focus on the three SoCalGas studies combined into an overall SoCalGas sample set for further analysis. This combined SoCalGas sample set then forms the foundation for the probabilistic-based emission factors associated with the Decision Tree groupings.

**Figure 11: Leak Rate Cumulative Fraction of Two Sample-Biased SoCalGas Studies.**

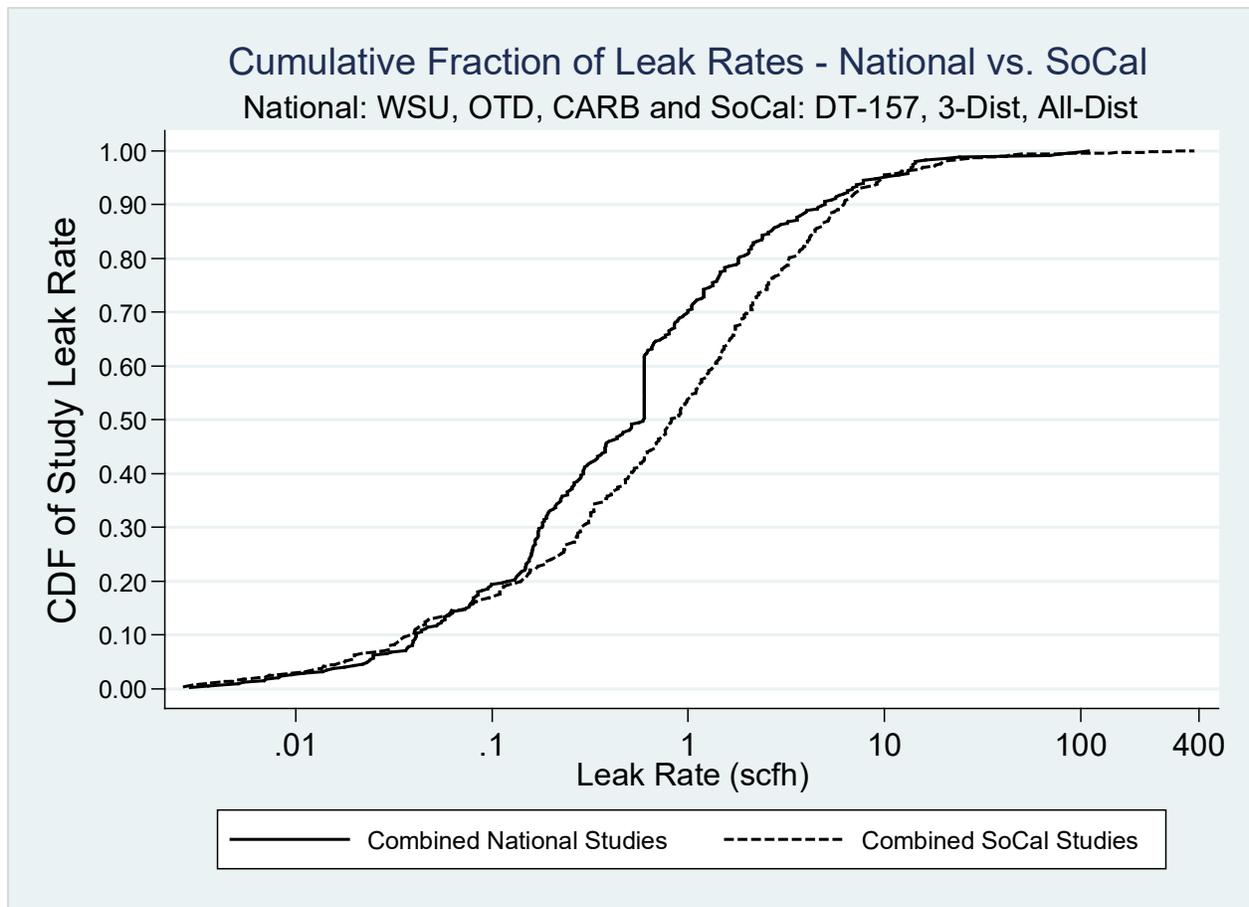


## Data Analysis using Cumulative Fraction Plots

A more quantitative review of the leak rate distributions can be performed through cumulative fraction plots. These are sometimes referred to as frequency diagrams or cumulative distribution functions when discussing theoretical distributions. The cumulative fraction plots of the combined national and SoCalGas studies are shown in Figure 12 below.

The distributions track on top of each other (on the ends) and diverge between approximately 0.20 and 0.90 fraction of the samples. The studies have effectively the same fractions (percentiles) of samples below 0.1 scfh and above 10 scfh. One could say they have the same fraction of "large" non-hazardous leaks, i.e. approximately 95% of the leaks are less than 10 scfh for both combined groups. Quantitative tests of sameness are presented later in this report.

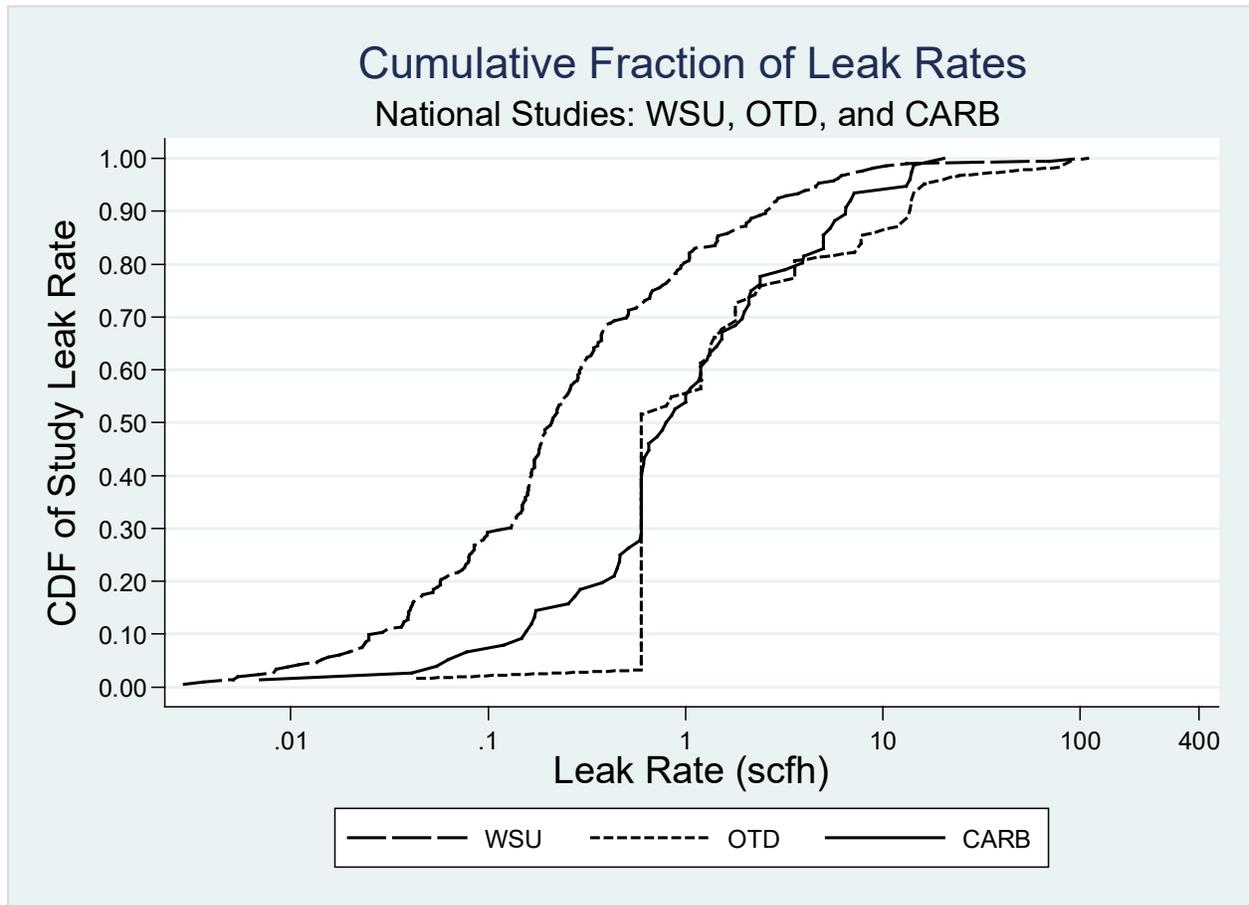
**Figure 12: Leak Rate Cumulative Fraction of Combined National and SoCalGas Studies.**



In Figure 13, the three national studies are plotted as cumulative fractions. The OTD and CARB studies are similar and intersect each other four separate times between 1 and 10 scfh, whereas the WSU study shows a lower overall distribution of leaks and is always to the left and above the

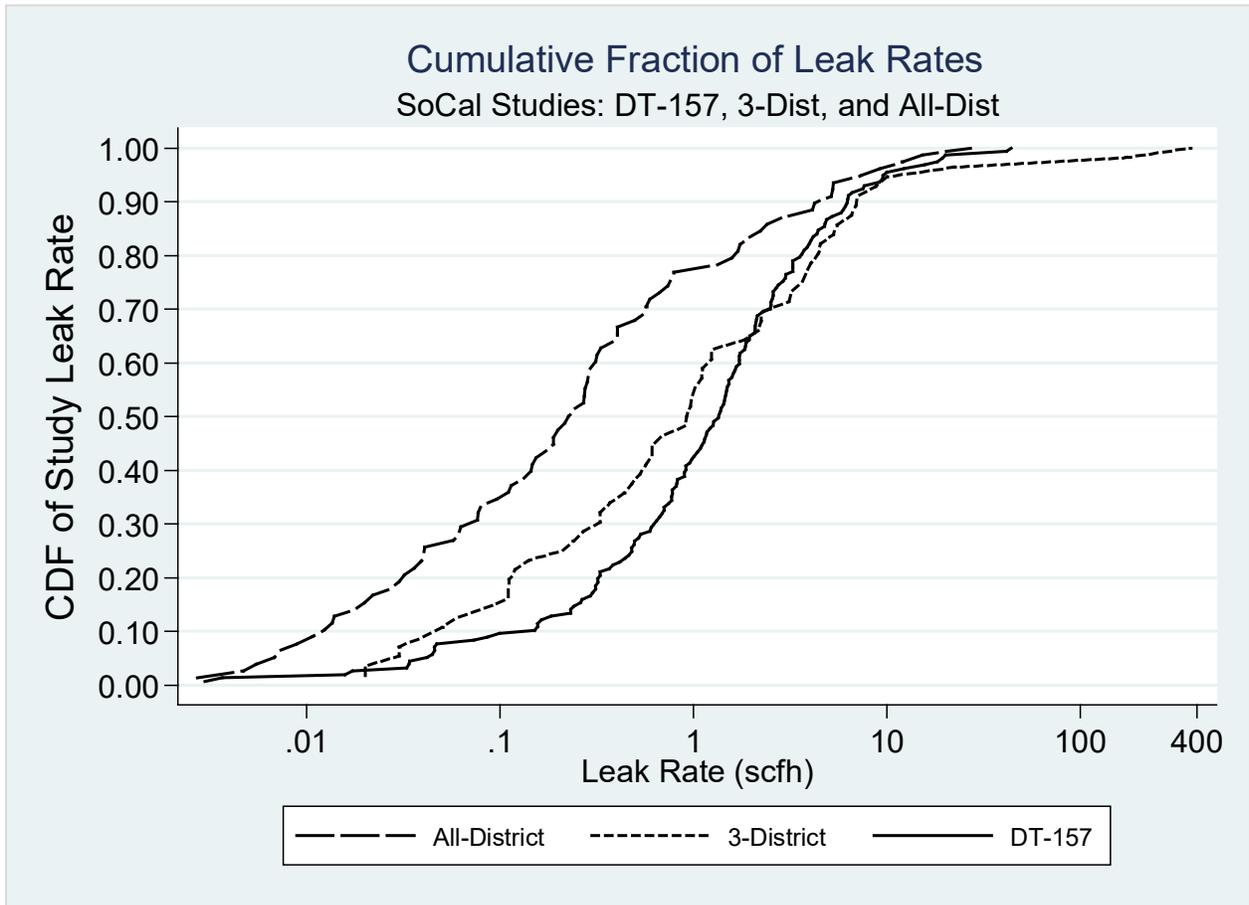
other two study cumulative fraction plots. This is probably due to the lower detection limit in the WSU measurements compared to the other studies.

Figure 13: Leak Rate Cumulative Fraction of National Studies.



Finally, in Figure 14, the three SoCalGas studies are plotted as cumulative fractions. The 3-District and DT-157 studies are similar and intersect each other three separate times between 1 and 10 scfh, whereas the All-District study shows a lower overall distribution of leaks and is always to the left and above the other two study cumulative fraction plots.

Figure 14: Leak Rate Cumulative Fraction of Three SoCalGas Studies.



### 5.3. Data Transformation

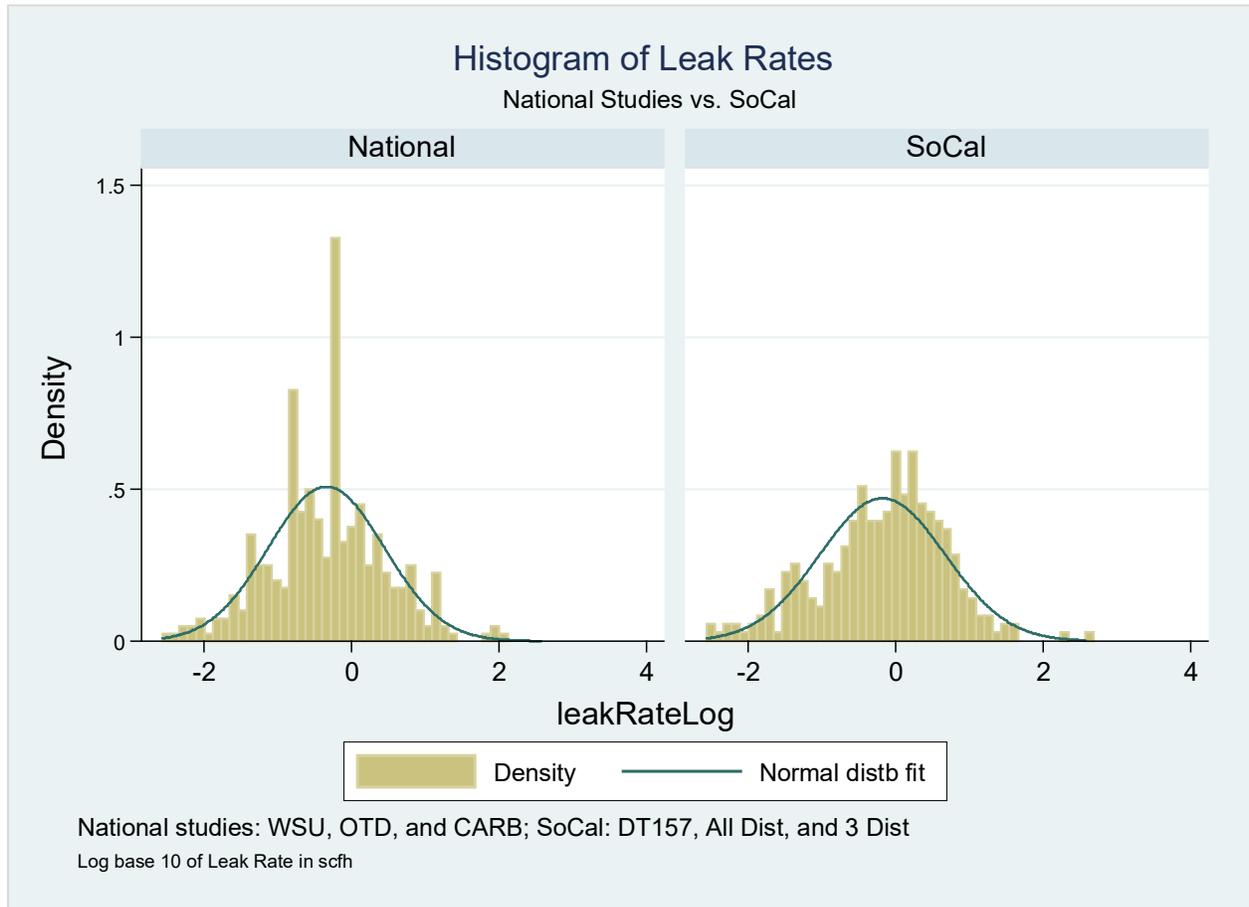
In prior sections, the national and SoCalGas data sets were presented, the data distributions were explained, and the studies with known (a priori) sampling bias (which were shown through the descriptive statistics) were removed.

Prior to running regression and probabilistic analysis to check for outliers and/or other issues with the data sets, this section analyzes the combined SoCalGas and combined national studies (as a baseline) to the leak rate distributions to determine if a transformation can be applied to make the rates normally distributed. The combined national studies are done for comparative purposes. Based off of the three analysis further described below, the log (base 10) of the leak rate is shown to be an appropriate transformation and will therefore be used in regression analysis, as well as for residual and diagnostic analysis.

#### Histogram of Log Transformed Leak Data

The national and SoCalGas combined sets were transformed to  $\log(10)$  of the leak rate and plotted as histogram distributions. These are shown side-by-side in Figure 15 below. A normal distribution fit is overlaid on the density plots and shows good agreement. However, a more quantitative measure is needed.

**Figure 15: Leak Rate Histogram of Log(10) of Combined National and SoCalGas Studies.**



### One-sample Kolmogorov-Smirnov Test for Lognormality

A Kolmogorov-Smirnov (KS) test for normality of the log(10) of leak rate was conducted on both the combined national studies, see Table 6, and for the combined SoCalGas studies, see Table 7. For the Kolmogorov-Smirnov (K-S) test, a combined K-S p-value of greater than 0.05 for a one sample test of normality indicates there is no support that the distribution is not normal distributed.

The K-S test is set up and reported in a manner opposite than most null hypotheses which are stated in a way that there is no relationship (i.e. the results are random) between two variables being studied. Both statistical tests result in supporting the possibility that the log transformed leak rate data is normally distributed as shown by their respective p-value being greater than 0.05.

**Table 6: Kolmogorov-Smirnov test of Combined National Studies for Normality.**

One-sample Kolmogorov-Smirnov test against theoretical distribution normal((leakRateLog+.3443504)/.7850414)		
Smaller group	D	P-value
-----		
leakRateLog:	0.0576	0.098
Cumulative:	-0.0622	0.067
Combined K-S:	0.0622	<b>0.133</b>

**Table 7: Kolmogorov-Smirnov test of Combined SoCalGas Studies for Normality.**

One-sample Kolmogorov-Smirnov test against theoretical distribution normal((leakRateLog+.175378)/.857963)		
Smaller group	D	P-value
-----		
leakRateLog:	0.0407	0.382
Cumulative:	-0.0574	0.147
Combined K-S:	0.0574	<b>0.292</b>

## Quantile-Normal Plot Analysis of Transformations

To visually check the transformation analysis, the quantile-normal plots of the log(10) transformed data were plotted for the combined national studies, see Figure 16, and for the combined SoCalGas studies, see Figure 17. Both plots show that leak rate data from the combined data sets appears log-normal distributed.

Figure 16: Quantile-Normal Plot of Log-normal Transformed Combined National Studies.

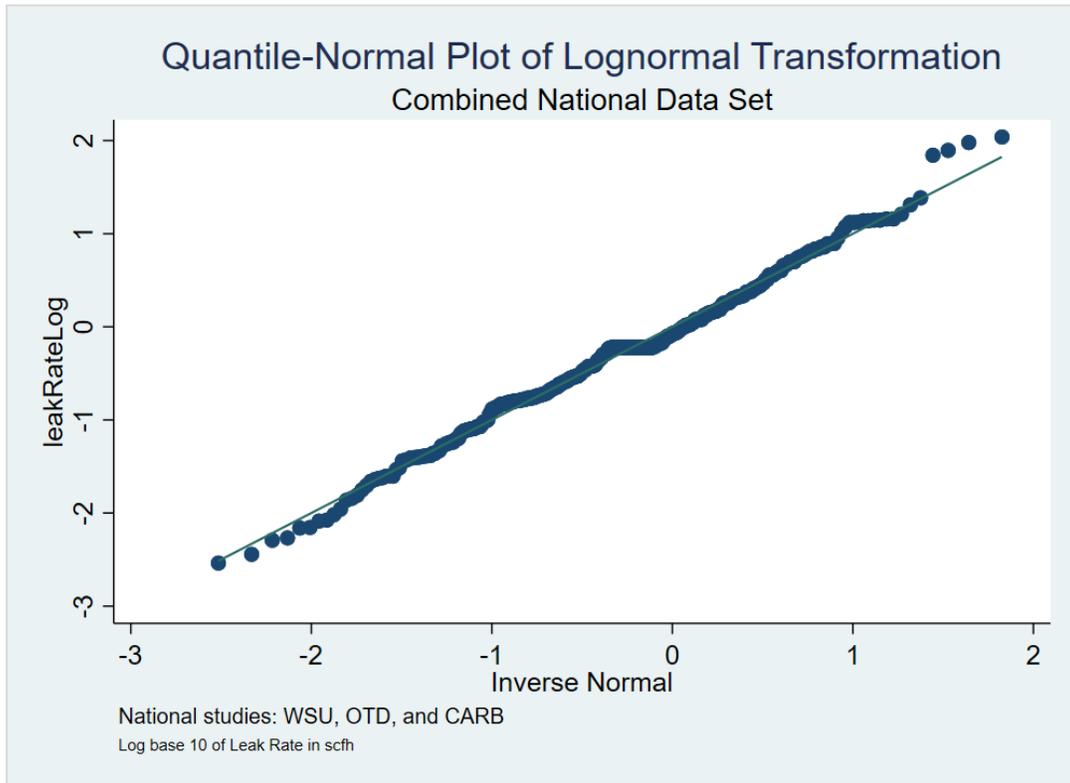
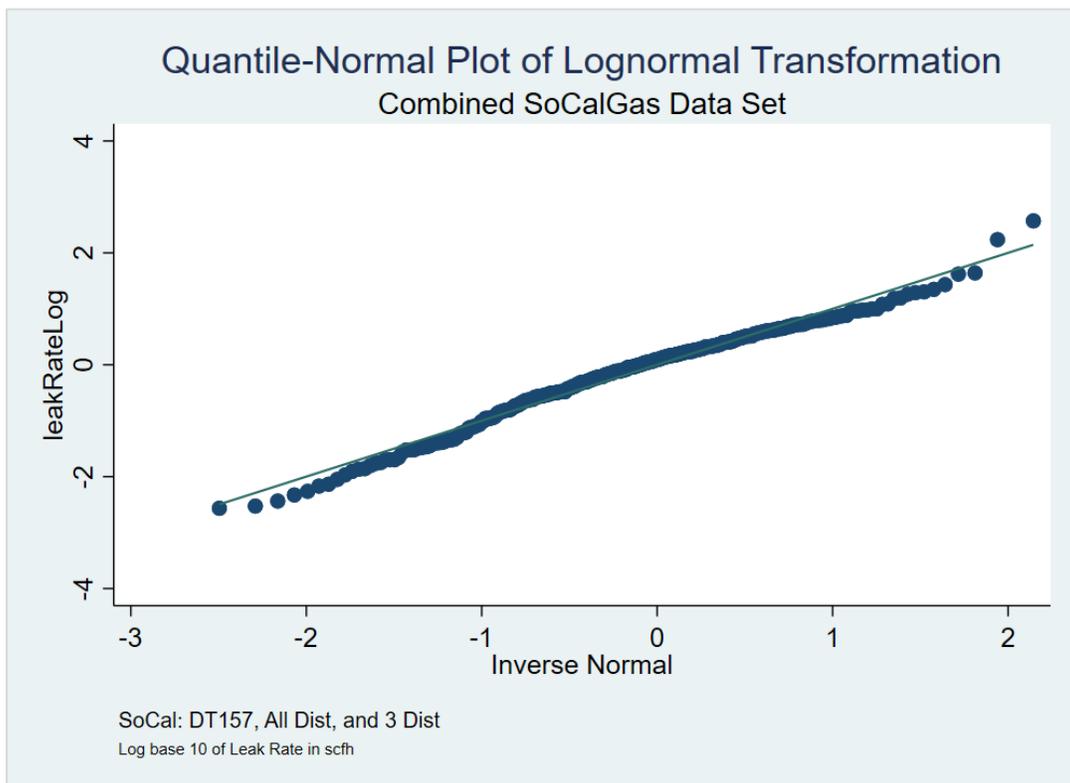


Figure 17: Quantile-Normal Plot of Log-normal Transformed Combined SoCalGas Studies.



## 5.4. Linear Regression Study Means Comparison

### Combined National vs. Combined SoCalGas Study Means Analysis

The ANOVA analysis and output table for their first use are explained below. For the rest of this report, key conclusions from the ANOVA will be discussed with a detailed table included in the corresponding section or in the Appendix as noted.

#### Analysis of Variance (ANOVA)

Using the log transformation for leak rate, an ANOVA was completed between the combined national and all SoCalGas studies. The results are in Table 8 below. ANOVA is used to determine any difference in a metric variable ( $\log(10)$  leak rate in this case), between two or more groups.

The upper part of the table shows the means for each group where one can observe the difference already noted in the descriptive statistics section of this report. However, note this is the mean of the  $\log(10)$  of each measurement which is not the same as the  $\log(10)$  of the mean of the leak measurements. The Prob > F or p-statistic is smaller (0.0099) than 0.05 and therefore indicates the result is *significant*, which means one would expect the same result (difference in means) if the entire *population* of the applicable studies were measured.

It is worth noting for completeness that the p-statistic for significance is derived from the F distribution - where F is the ratio of the two variance estimates (MS, or mean squares) listed in the ANOVA table for “between” vs. “within groups” respectively. The MS are calculated by the ratio of the sum of the squares (of deviation) for each source to the degrees of freedom.

Therefore, one can conclude that at least two groups exhibit a statistically significant difference in their means. Another way to interpret this would be to say that one would expect less than 1 time in 100 that these results would be obtained if there were no difference between the national and SoCalGas combined studies.

*Bartlett's* test for equal variances indicates a non-significant result which is good; otherwise, one would have to conclude that the variance between the groups were unequal. Also note the standard deviations between the two combined groups is close (0.79 and 0.85 in round numbers), and the frequency are both large and similar (350 and 309).

Finally, *Bonferroni* compares between all possible groups. In this case, it is simply a one-to-one comparison, so it is the same as the ANOVA. Additionally, the difference between  $\log(10)$  means is shown with the same p-statistic as the overall ANOVA, but to only three vs. four significant figures (0.010).

**Table 8: ANOVA of Combined National vs. Combined SoCalGas Leak Rate Means.**

Summary of Log of Leak Rate (scfh)					
study Scale	Mean	Std. Dev.	Freq.		
National	-.3443504	.78504136	350		
SoCal	-.17955489	.84894047	309		
Total	-.26707906	.81914558	659		
Analysis of Variance					
Source	SS	df	MS	F	Prob > F
Between groups	4.45688946	1	4.45688946	<b>6.70</b>	<b>0.0099</b>
Within groups	437.060768	657	.665237089		
Total	441.517657	658	.670999479		
Bartlett's test for equal variances: chi2(1) = 2.0052 Prob>chi2 = 0.157 Comparison of Log of Leak Rate (scfh) by studyScale (Bonferroni)					
Row Mean-					
Col Mean	National				
	SoCal	.164796			
		0.010			

The ANOVA includes several assumptions worth noting:

- The outcome variable is quantitative - true in this case.
- The errors or residuals are normally distributed. This could be problematic with small sample sets. It will be shown that the residuals are normally distributed in a later section of this report.
- The observations represent a random sample of the population. The two sample-biased studies which also contained ten or less samples were removed - which would be a problem to determine the above residual requirement as well, i.e. to quantify normality.
- The errors are independent. A good assumption in this study.
- The variance of each group is equal. This is the aforementioned Bartlett test which confirmed this assumption.

The ANOVA across the individual studies was run and confirmed that there was a difference between individual studies (not just the difference between the combined studies noted above). The results showed a  $F(7, 651) = 17.79$   $p < 0.001$  meaning there is statistical difference between the study means. For those interested, the ANOVA table by study and a full pairwise comparison by study-to-study is presented in Appendix B. Instead of going into the individual study comparisons to one another with the ANOVA results, these will be discussed in the *regression*

section below, which also allows the use of diagnostic tests to pinpoint outliers or extreme values in the data.

## Linear Regression

### Overview

The dependent (and continuous) variable for the linear regressions (LR) is the log(10) leak rate. The independent (and categorical) variable for the regression is the emission study for each leak observation (sample).

The linear regression results for the log(10) of leak rate by study is shown in Table 9 below. The first section of the regression output shows the model and residual (these are termed sources) sum of the squares of deviation, degrees of freedom, and mean square for the sources. The right side of the table provides the total number of sample observations, the F test score (7, 651) for the model and associated residual of 17.79, and the associated p-statistic which is less than 0.001. The number to focus on is the p-statistic which was explained earlier in the ANOVA section of the report. This means a significant model is observed, and the R-squared and adjusted R-squared values show a moderate level of accounting for variance.

The bottom section lists the regression coefficients (for all independent categorical variables) with the mean log(10) leak rates by study and associated error terms. These coefficients are used to calculate the t-value and the p-statistic. Again, placing focus on the p-statistic, it can be seen that for the non-biased studies, the differences in means for the All District Study and the WSU studies cannot be explained by random variation in the samples, i.e. the differences are significant since their p-statistic is less than 0.05 - as was discussed in the earlier sections of this report. Further, the All District Study and the WSU Study are statistically similar as noted in their pairwise large p-significance level of 0.678 (see Appendix B for regression details).

**Table 9: Linear Regression of Individual National and SoCalGas Study Leak Rate Means.**

Source	SS	df	MS	Number of obs	=	659
Model	70.8808391	7	10.1258342	F(7, 651)	=	17.79
Residual	370.636818	651	.56933459	Prob > F	=	0.0000
				R-squared	=	0.1605
				Adj R-squared	=	0.1515
Total	441.517657	658	.670999479	Root MSE	=	.75454
-----						
leakRateLog	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
-----						
study						
3DisPilot	0	(base)				
3DisPilotLowSpec	-.5881559	.2851903	-2.06	0.040	-1.14816	-.028152
AllDisLIRP	.1388059	.259037	0.54	0.592	-.3698429	.6474547
AllDisPilot	-.5897834	.1321584	-4.46	0.000	-.8492916	-.3302752
DT157Pilot	.0843352	.1174437	0.72	0.473	-.146279	.3149494
Natl_CARB_GTI	.0444857	.1328832	0.33	0.738	-.2164457	.305417
Natl_OTD_GTI	.230709	.1391025	1.66	0.098	-.0424347	.5038527
Natl_WSU_EDF	-.5482558	.1133677	-4.84	0.000	-.7708662	-.3256454
_cons	-.0627922	.10083	-0.62	0.534	-.2607835	.135199

### Three SoCalGas Study Means Analysis

In this section, the ANOVA technique is applied to the three SoCalGas studies that will provide the basis for the subsequent emission factor calculations.

#### Analysis of Variance (ANOVA)

In Table 10 below, there exists similar deviations between the studies and large sample sizes. The F test score (2, 288) = 18.72 and the p-statistic is less than 0.001. This result demonstrates evidence that statistical differences are present in the means of the log(10) leak rates between the studies in the population. The difference can be seen in the pairwise combinations at the bottom of the table. The AllDisPilot is different from the other two which are similar based on the p-statistic values.

**Table 10: ANOVA of Leak Rate Means for Three SoCalGas Studies.**

Summary of Log of Leak Rate (scfh)					
study	Mean	Std. Dev.	Freq.		
3DisPilot	-.06279222	.89898228	56		
AllDisPil	-.65257563	.92139692	78		
DT157Pilo	.02154301	.71202596	157		
Total	-.17537804	.85796296	291		

Analysis of Variance					
Source	SS	df	MS	F	Prob > F
Between groups	24.5599302	2	12.2799651	18.72	0.0000
Within groups	188.909199	288	.655934719		
Total	213.469129	290	.736100446		

Bartlett's test for equal variances:  $\chi^2(2) = 8.7913$  Prob> $\chi^2 = 0.012$

Comparison of Log of Leak Rate (scfh) by study (Bonferroni)			
Row Mean- Col Mean	3DisPilo	AllDisPi	
AllDisPi	-.589783 0.000		
DT157Pil	.084335 1.000	.674119 0.000	

## Linear Regression

The dependent (and continuous) variable for the linear regressions is the log(10) leak rate. The independent (and categorical) variable for the regression is the emission study for each leak observation (sample).

The linear regression of the SoCalGas studies is shown in Table 11 below and further shows that the DT-157 and 3-District study means for log(10) leak rate are statistically the same, and that the All District study mean has a statistically significant difference (a lower value) from both of the other two studies. This was explained in an earlier section of this report when looking at the raw data and descriptive statistics and is not an anomaly that should be filtered out.

**Table 11: LR and PW Comparison of Leak Rate Means for Three SoCalGas Studies.**

Source	SS	df	MS	Number of obs	=	291
				F(2, 288)	=	18.72
Model	24.5599302	2	12.2799651	Prob > F	=	0.0000
Residual	188.909199	288	.655934719	R-squared	=	0.1151
				Adj R-squared	=	0.1089
Total	213.469129	290	.736100446	Root MSE	=	.8099
-----						
leakRateLog	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
-----						
study						
3DisPilot	0	(base)				
AllDisPilot	-.5897834	.1418539	-4.16	0.000	-.8689853	-.3105815
DT157Pilot	.0843352	.1260597	0.67	0.504	-.1637799	.3324504
_cons	-.0627922	.1082272	-0.58	0.562	-.2758087	.1502243
-----						
Pairwise comparisons of marginal linear predictions						
-----						
	Contrast	Std. Err.	Unadjusted t	P> t	[95% Conf. Interval]	
-----						
study						
AllDisPilot vs 3DisPilot	-.5897834	.1418539	-4.16	0.000	-.8689853	-.3105815
DT157Pilot vs 3DisPilot	.0843352	.1260597	0.67	0.504	-.1637799	.3324504
DT157Pilot vs AllDisPilot	.6741186	.1121933	6.01	0.000	.4532957	.8949415

## Linear Regression Residual Analysis and Regression Diagnostics

A series of diagnostics were analyzed for the SoCalGas study regressions to confirm regression assumptions and to look for influential, outlier, or extreme values requiring further review and/or explanation or exclusion. All data points were retained after this detailed review of the diagnostics and residuals, including:

- Exogeneity
- Random Sampling
- Linearity in Parameters
- Multicollinearity
- Heteroscedasticity and Normal Distribution of Residuals
- Influential Observations - DFBETA
- Influential Observations - Cook's Distance
- Influential Observations - Leverage

The details and plots of all the diagnostics and residuals are presented in Appendix B.

## 5.5. Bayesian Monte Carlo Markov Chain (MCMC) Regression

A non-parametric Bayesian Monte Carlo Markov Chain (MCMC) [25] *regression* analysis was conducted. Two variations of random *sampling* were used, the random walk Metropolis-Hasting (MHS) [26, 27] method as well as the more robust Gibbs (GS) [28] method. In both cases, 35,000 iterations were used with a 5,000 iteration burn-in run, resulting in an incorporated Monte Carlo sample size of 30,000. In both cases, the prior distribution for the  $\log(10)$  leak rate distribution was "uniformed", i.e. a flat/uniform prior. The sigma prior was assumed as a conservative gamma function.

The results of these Bayesian-based regressions showed very similar results in output to the standard regression outputs already discussed. However, this is for convenience of comparison, since the methods are completely different, and this analysis uses Bayesian linear regression. This does not come as a surprise, since the regression assumptions were met, and the dependent variable ( $\log(10)$  of leak rate) was normally distributed.

The detailed results of the Bayesian MCMC(MHS) and MCMC(GS) are presented in Appendix B.

## 5.6. Sensitivity of Leak Rate to Geographic District and Year of Detection Analysis

The sections below describe two sensitivity studies that were completed to determine if the SoCalGas studies demonstrated sensitivity in the log(10) leak rate values to *geographic* district (location of leak) and/or the *year* the leak was originally detected. The geographic districts and year of leak detection are not listed in the Appendix table. As shown below, the leak rate was not sensitive to either the geographic districts or the year the leak was originally detected.

### Geographic District of Leak

The ANOVA analysis for geographic sensitivity is shown in Table 12 below. The p-statistic of 0.1087 shows that the difference between groups (districts) is not statistically significant; hence, there is no evidence of a difference in the log(10) leak rate values between different geographic districts. Note that for the analysis, several districts had to be dropped out of this particular ANOVA analysis due to only one sample with leak rates present. These were not dropped out of the overall study.

**Table 12: ANOVA of Leak Rate Means Across Districts for Three SoCalGas Studies.**

Source	Analysis of Variance			F	Prob > F
	SS	df	MS		
Between groups	69.7037636	83	.839804381	1.24	0.1087
Within groups	152.271814	225	.676763619		
Total	221.975578	308	.720699928		

Bartlett's test for equal variances:  $\chi^2(47) = 69.4664$  Prob> $\chi^2 = 0.018$

note: Bartlett's test performed on cells with positive variance:  
36 single-observation cells not used

### Year of Leak Detection

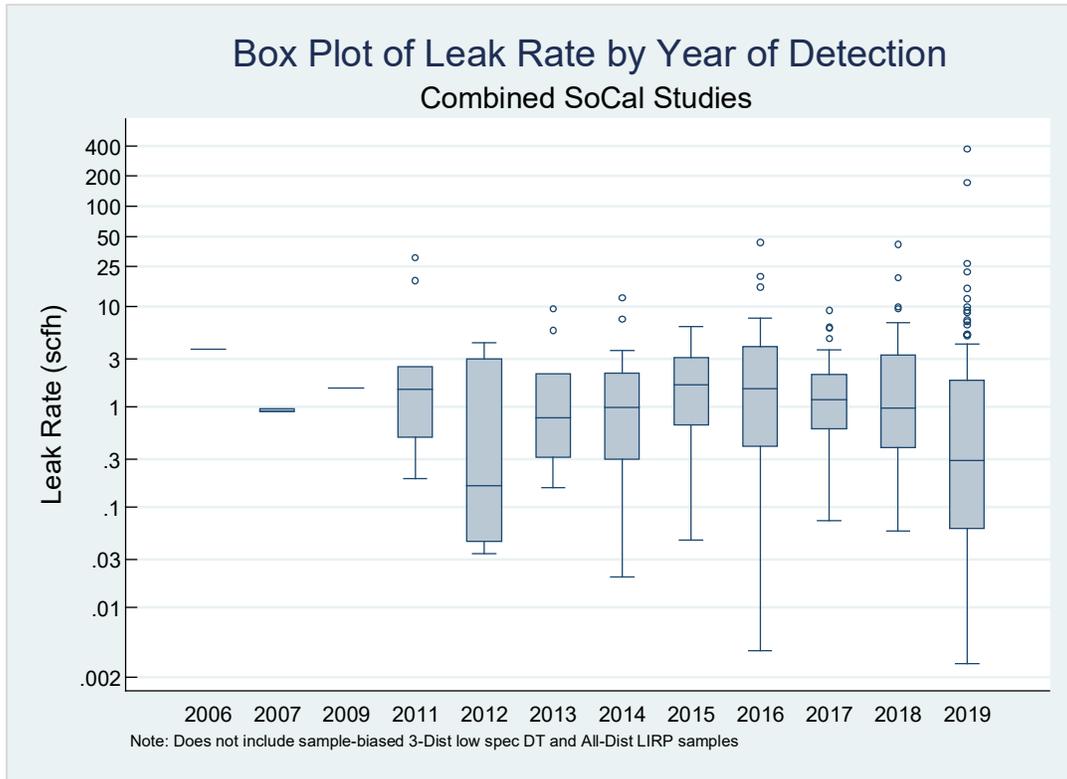
To get an understanding for the distribution of leak rates by year of detection, two plots were generated. The first is a box plot of leak rate by year detected (Figure 18). A scatter plot of the same data (Figure 19) shows the median, mean, and maximum leak rates detected by year.

Datasets from 2006, 2007, 2009, and 2012 were removed from this regression analysis since they only have 1, 3, 1, and 5 observations respectively. They were not removed from the study overall.

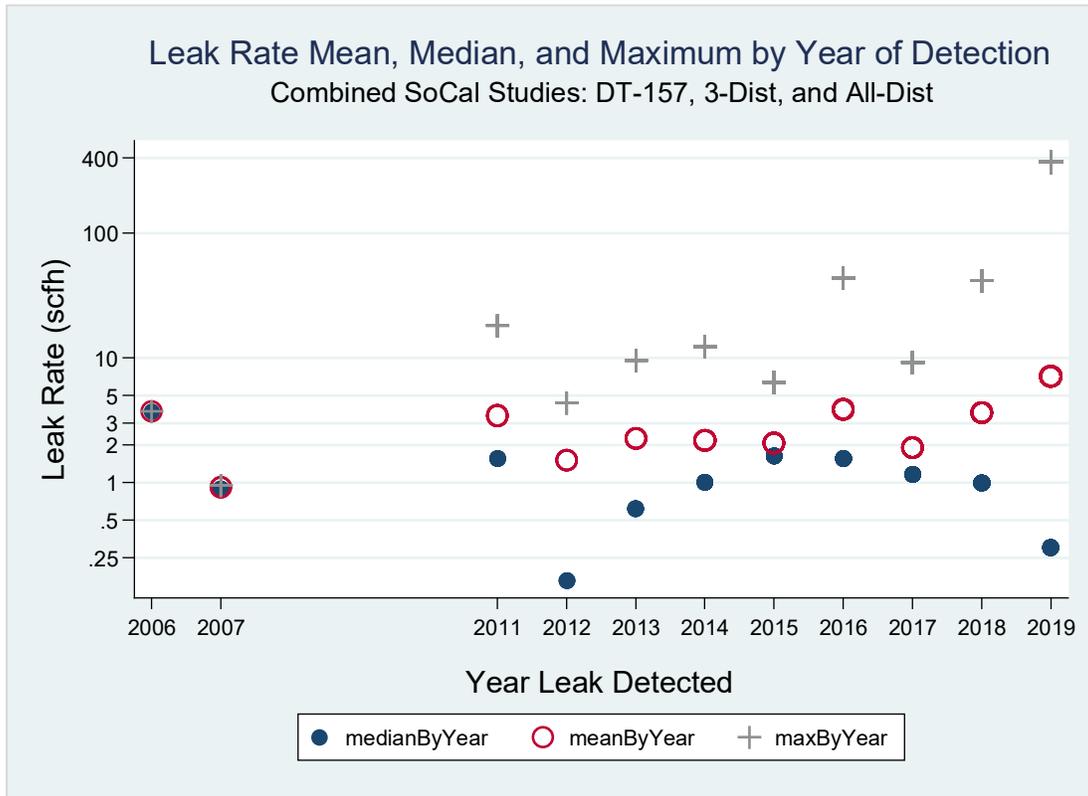
From the plots, one can see that 2019 has the most data, accounting for over 1/3 of all observations. It also has the smallest and highest single values of leak rates and the lowest median.

The medians of leak rate by year detected are very consistent, all between 0.5 and 2.0 scfh, except for 2019 which is below 0.50 scfh. The highest median leak rate occurred in 2015 after which the median leaks continue to decrease by year.

**Figure 18: Leak Rate Box Plots by Year Leak Detected for Three Combined SoCalGas Studies.**



**Figure 19: Leak Rate Median, Mean, and Maximum by Year SoCalGas Leaks Detected.**



The ANOVA analysis is shown in Table 13 and the pairwise comparison in Table 14. Both of these results, show a quantitative measure of sameness or lack thereof.

In summary, the log(10) leak rate is insensitive to geographic district or year of detection based on the field data analyzed to date.

**Table 13: ANOVA of Leak Rate Means Across Year Detected for Three SoCalGas Studies.**

Summary of Log of Leak Rate (scfh)			
Year Leak Detected	Mean	Std. Dev.	Freq.
2006	.57175308	0	1
2007	-.03804038	.01594704	3
2009	.18752073	0	1
2011	.18185936	.67890395	11
2012	-.49550171	.99638628	5
2013	-.04097877	.59161147	10
2014	-.13138904	.74362153	19
2015	.08858402	.54923189	32
2016	.0052503	.84882588	39
2017	.00980475	.53934181	29
2018	.02103485	.67407974	38
2019	-.45726365	1.0032278	111

-----					
Total	-.16512119	.84288247		299	
Analysis of Variance					
Source	SS	df	MS	F	Prob > F
-----					
Between groups	17.6311516	11	1.60283196	2.37	0.0081
Within groups	194.083206	287	.676248105		
-----					
Total	211.714358	298	.710450864		
Bartlett's test for equal variances: chi2(9) = 43.0070 Prob>chi2 = 0.000					
note: Bartlett's test performed on cells with positive variance:					
2 single-observation cells not used					

**Table 14: PW Comparison of Leak Rate Means by Year Detected for Three SoCalGas Studies.**

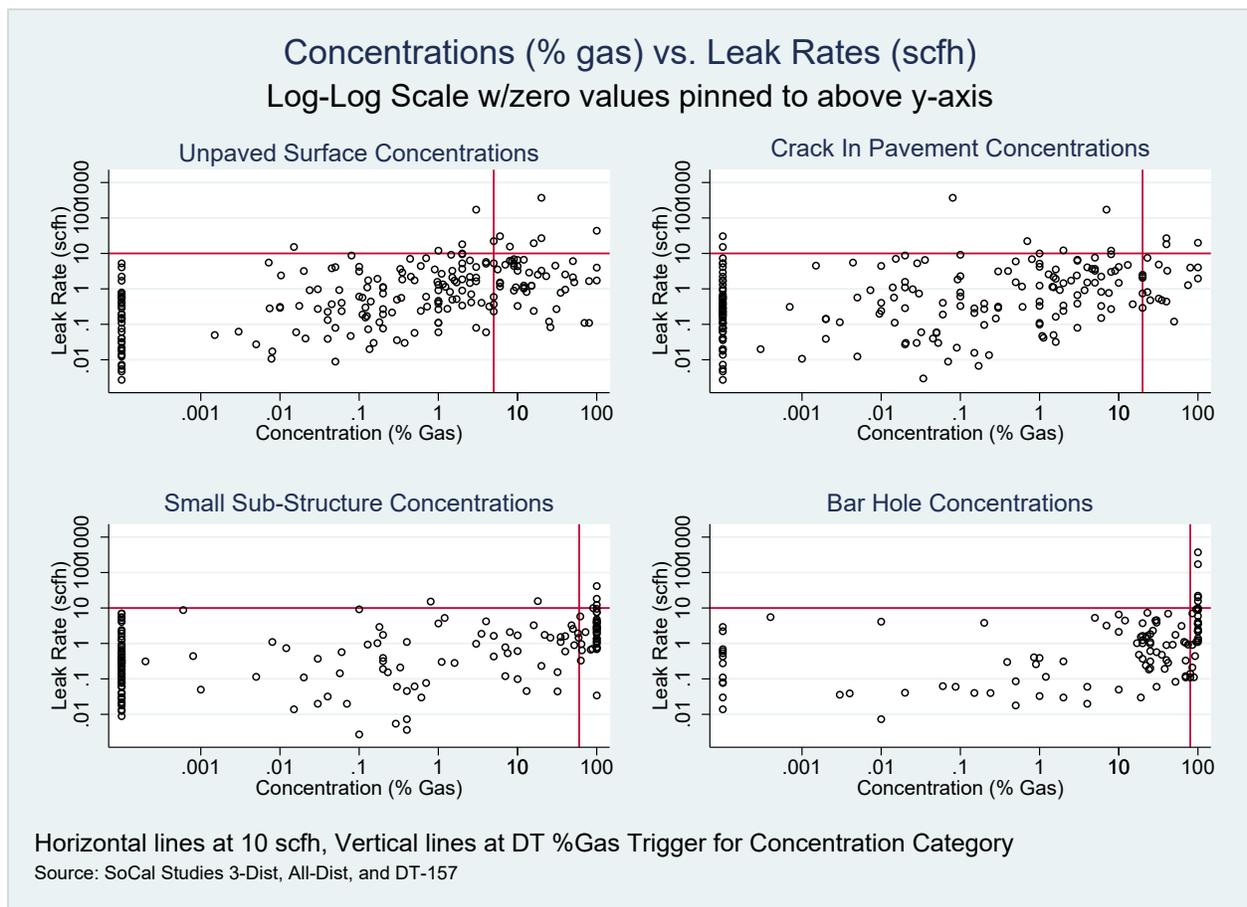
Comparison of Log of Leak Rate (scfh) by yearDetected (Bonferroni)						
Row Mean- Col Mean	2006	2007	2009	2011	2012	2013
2007	-.609793 1.000					
2009	-.384232 1.000	.225561 1.000				
2011	-.389894 1.000	.2199 1.000	-.005661 1.000			
2012	-1.06725 1.000	-.457461 1.000	-.683022 1.000	-.677361 1.000		
2013	-.612732 1.000	-.002938 1.000	-.2285 1.000	-.222838 1.000	.454523 1.000	
2014	-.703142 1.000	-.093349 1.000	-.31891 1.000	-.313248 1.000	.364113 1.000	-.09041 1.000
2015	-.483169 1.000	.126624 1.000	-.098937 1.000	-.093275 1.000	.584086 1.000	.129563 1.000
2016	-.566503 1.000	.043291 1.000	-.18227 1.000	-.176609 1.000	.500752 1.000	.046229 1.000
2017	-.561948 1.000	.047845 1.000	-.177716 1.000	-.172055 1.000	.505306 1.000	.050784 1.000
2018	-.550718 1.000	.059075 1.000	-.166486 1.000	-.160825 1.000	.516537 1.000	.062014 1.000
2019	-1.02902 1.000	-.419223 1.000	-.644784 1.000	-.639123 0.959	.038238 1.000	-.416285 1.000
Row Mean- Col Mean	2014	2015	2016	2017	2018	
2015	.219973 1.000					
2016	.136639 1.000	-.083334 1.000				
2017	.141194 1.000	-.078779 1.000	.004554 1.000			
2018	.152424 1.000	-.067549 1.000	.015785 1.000	.01123 1.000		
2019	-.325875 1.000	-.545848 0.070	-.462514 0.181	-.467068 0.452	-.478298 0.143	

## 5.7. Concentration vs. Leak Rate Analysis

### General Trends

The four leak concentration measurement categories are plotted separately, but side-by-side in Figure 20. The Decision Tree 10 scfh leak rate value that separates “Not Large” from “Large” non-hazardous leak levels is plotted as a horizontal line, and each of the concentration levels that trigger a positive Decision Tree categorization are plotted as a vertical line at 80%, 20%, 60%, and 5% gas respectively. The zero values for concentration were set to 0.0001% gas which is less than the minimum value, appearing on the left-hand side of the plots.

Figure 20: Separate Plots of Leak Concentrations vs. Rates by DT Category.



One can see the general upward trend between the methane concentration measurements with increasing leakage flow rate. The scatterplots of maximum surface concentration vs. leak rate do support the thresholds. The upper left quadrants of the leak rate vs. concentration plots are the areas of false negatives for the DT process if each surface category was evaluated individually. One can visually see the very low number of false negatives (the measure of importance for the entire DT program) in these zones even when evaluated individually. In practice however,

because the DT is an 'OR' gated process (i.e., the DT will be triggered if *any* one of the four category thresholds are met) this substantially reduces the number of false negatives in actuality from what is shown individually. These observations are also fully supported by the Bayesian probabilistic analysis later described in this report.

A regression analysis was done for general trend review. The regression analysis was *not* used for subsequent quantitative calculations related to the Decision Tree predictive capability, which is covered in the next section using Bayesian probabilistic analysis. The analysis demonstrated that any individual concentration measurement in any category is *not* a good predictor of leakage flow rate. The regression analysis is presented in Appendix B of this report as supplemental information.

## 5.8. Decision Tree Leak Prediction Quantitative Performance

Next, the Decision Tree performance measures were developed as related to its ability to properly predict a greater or equal to 10 scfh leak rate given any combination from one to all four measurement concentration measurement categories for any individual leak.

### Leak Rate Statistics of Empirical Data by Decision Tree Groupings

#### Overall Results from Empirical Data

The combined SoCalGas study mean leak flow rate is presented again for the 291 samples in Table 15 with the associated 95% confidence intervals. As noted earlier in the report, the assumption of normality for leak rate data is violated, since we have a highly skewed distribution. The mean values and associated 95% confidence intervals presented in Table 15, Table 20, and Table 23 below are from the non-parametric bootstrap analysis conducted in Section 5.10 of the report.

**Table 15: Combined Bootstrap Mean and C.I. for Three SoCalGas Studies (Baseline).**

SoCal Studies	Obs	Mean	[95% Conf. Interval]	
Combined	291	4.303	1.635	12.013

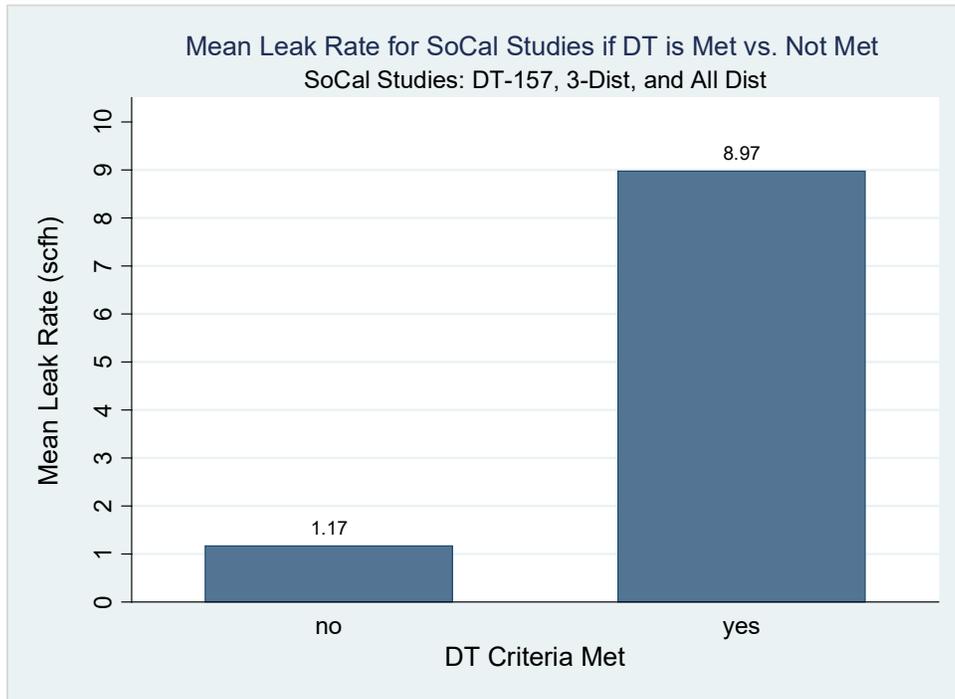
For comparison, Table 16 presents the mean, minimum, and maximum in two groups: when the DT was not met and when the DT was met from the concentration value(s).

Excellent separation of the mean values is observed which are 1.168 scfh and 8.973 scfh respectively. This is plotted in Figure 21, and the median and 5<sup>th</sup> and 95<sup>th</sup> percentiles are shown in Table 17.

**Table 16: Leak Rate Mean, Min., and Max. of by DT Grouping.**

DT Met	N(count)	mean(scfh)	min(scfh)	max(scfh)
no	174	1.168	0.003	15.252
yes	117	8.973	0.034	373.000
<b>Total</b>	<b>291</b>	<b>4.306</b>	<b>0.003</b>	<b>373.000</b>

**Figure 21: Mean Leak Rate for DT Categories.**



**Table 17: Leak Rate Median, 5%, and 95% Percentiles by DT Grouping.**

DT Met	N(count)	p5(scfh)	med(scfh)	p95(scfh)
no	174	0.011	0.376	5.156
yes	117	0.120	2.125	22.290
<b>Total</b>	<b>291</b>	<b>0.018</b>	<b>0.824</b>	<b>9.990</b>

## Leak Rate Statistics of Empirical Data by Confirmed Leak Rate Groupings

In this section, *actual* leak rates are grouped into those less than 10 scfh and those greater than or equal to 10 scfh and then by whether or not the DT criteria was not met and was met.

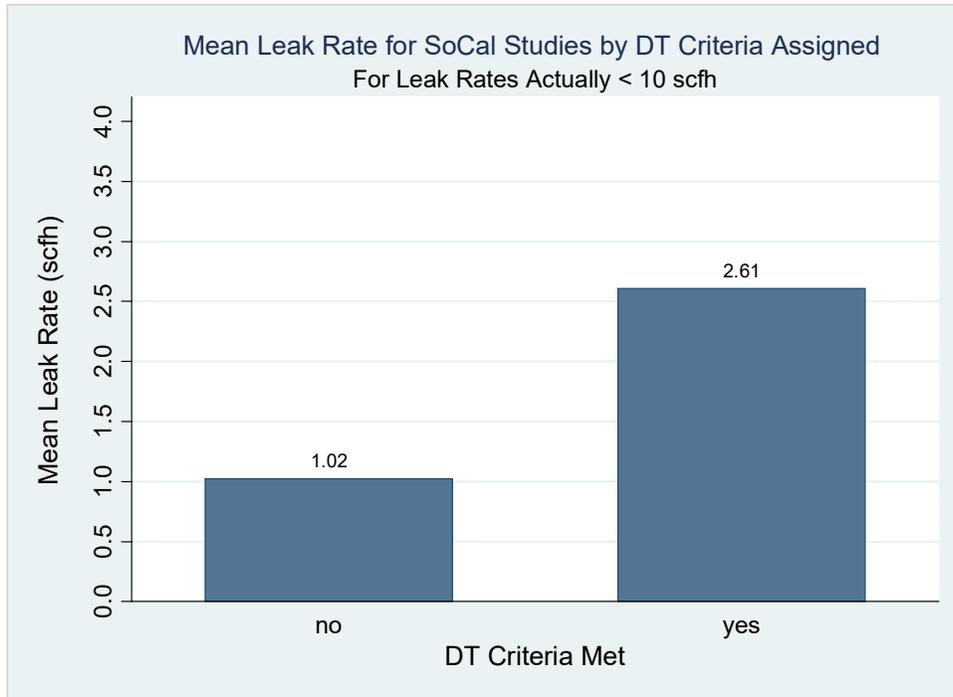
### Actual Leak Rate less than 10 scfh and True and False Negatives

When the DT is not met (i.e., predicting leak rate will be less than 10 scfh), we will refer to this as a "negative" prediction by the DT, and it is either true or false as determined by the subsequent leak rate measurement. In Table 18, the mean leak rate for these DT-related true and false negative situations are determined to be 1.023 scfh and 2.608 scfh respectively. These are plotted in Figure 22. The median and 5<sup>th</sup> and 95<sup>th</sup> percentiles are shown in Table 19.

**Table 18: Leak Rate Mean, Min., and Max. of Confirmed <10 scfh by DT Grouping.**

DT Met	N(count)	mean(scfh)	min(scfh)	max(scfh)
no	172	1.023	0.003	9.192
yes	105	2.608	0.034	9.990
<b>Total</b>	<b>277</b>	<b>1.624</b>	<b>0.003</b>	<b>9.990</b>

**Figure 22: Mean Leak Rate for DT Categories when Actually < 10 scfh.**



**Table 19: Leak Rate Med., 5%, and 95% Percentiles of Confirmed <10 scfh by DT Grouping.**

dt_met	N(count)	p5(scfh)	med(scfh)	p95(scfh)
no	172	0.011	0.350	4.434
yes	105	0.120	1.663	7.615
<b>Total</b>	<b>277</b>	<b>0.017</b>	<b>0.769</b>	<b>6.348</b>

The bootstrap mean leak rate of 1.623 scfh and confidence interval for the *actual* negatives, i.e. leaks with a flow rate less than 10 scfh, are listed in Table 20.

**Table 20: Bootstrap Leak Rate Mean and C.I. for Confirmed < 10 scfh (Actual Negatives).**

SoCal Studies	Obs	Mean	[95% Conf. Interval]
<10 scfh	277	1.623	1.141 2.153

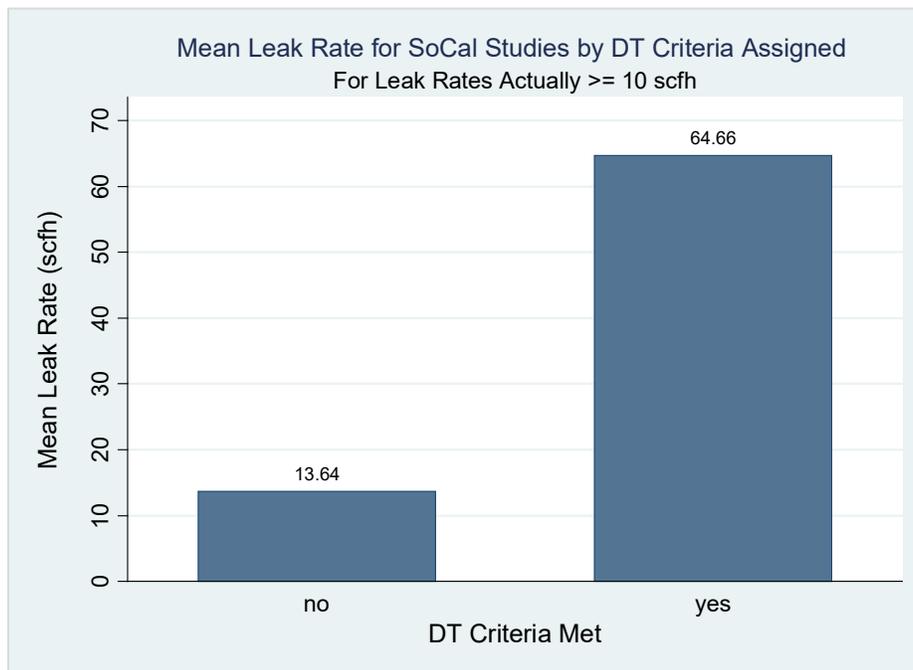
### Actual Leak Rate greater than or equal to 10 scfh

When the DT is met (predicting leak rate will be greater than or equal to 10 scfh), we will refer to this as a "positive" prediction by the DT, and it is either true or false as determined by the subsequent leak rate measurement. In Table 21, the mean leak rate for these DT-related true and false positive situations are shown. These are plotted in Figure 23. The median and 5<sup>th</sup> and 95<sup>th</sup> percentiles are shown in Table 22.

**Table 21: Leak Rate Mean, Min., and Max. of Confirmed  $\geq 10$  scfh by DT Grouping.**

DT Met	N(count)	mean(scfh)	min(scfh)	max(scfh)
no	2	13.638	12.024	15.252
yes	12	64.662	10.000	373.000
<b>Total</b>	<b>14</b>	<b>57.373</b>	<b>10.000</b>	<b>373.000</b>

**Figure 23: Mean Leak Rate for DT Categories when Actually  $\geq 10$  scfh.**



**Table 22: Leak Rate Med., 5%, and 95% Percentiles of Confirmed  $\geq 10$  scfh by DT Grouping.**

dt_met	N(count)	p5(scfh)	med(scfh)	p95(scfh)
no	2	12.024	13.638	15.252
yes	12	10.000	21.189	373.000
<b>Total</b>	<b>14</b>	<b>10.000</b>	<b>19.779</b>	<b>373.000</b>

The bootstrap mean leak rate of 57.667 and confidence interval for the *actual* positives, i.e. leaks with a flow rate greater than or equal to 10 scfh, are listed in Table 23.

**Table 23: Bootstrap Leak Rate Mean and C.I. for Confirmed  $\geq 10$  scfh (Actual Positives).**

<b>SoCal Studies</b>	<b>Obs</b>	<b>Mean</b>	<b>[95% Conf. Interval]</b>	
$\geq 10$	14	57.667	14.230	194.943

## 5.9. Bayesian Probabilistic Decision Tree Error-Type Analysis

A non-parametric Bayesian probabilistic analysis [15-17, 29] was conducted on the Decision Tree predictive power. The output includes the expected fraction (or percent) of sites that have true/false negative/positive outcomes. The Bayesian proportional analysis provides the most likely value of the errors in a coherent manner but also provides the upper and lower prediction limits around these values.

The results of the analysis are presented in Table 24 to Table 26 below, then plotted in Figure 24 and Figure 25. Table 24 and Figure 24 show the errors if one did not have prior knowledge of the leak concentration levels required for DT categorization as a likely large versus not large non-hazardous leak. In other words, it provides the likelihoods of any leak being in one of the four categories when concentration measurements were not available to input into the DT model. Normally, this would not be the case, since one will typically start with the informed knowledge of the DT category being met (positive) or not met (negative); however, it does provide a way to estimate the likely leak rate even without a concentration measurement available based on pure probability analysis.

Table 25 and Table 26 are of interest for this study since they provide the errors of a DT positive classification being true or false (Type I error) or of a DT negative classification being true or false (Type II error). These are plotted in Figure 25.

The DT has a low expected Type II error (false negative) of 1.1% and high, but conservative from an emissions standpoint, Type I error (false positive) of 89.7%. The lower and upper prediction (credible) limits are also tight, exhibiting a strong degree of belief and relatively low level of uncertainty.

### Joint False/True Positive (Type I) and Negative (Type II) Errors

**Table 24: Type I & II Uninformed (DT Cat. Unknown) Errors with 5% and 95% Pred. Limits.**

Error Type	Count	LPL%	MLV%	UPL%
False Neg	2	0.281	0.687	2.140
False Pos	105	31.621	36.082	40.841
True Neg	172	54.291	59.107	63.727
True Pos	12	2.653	4.124	6.573

### Independent False/True Positive Error Type I

**Table 25: Type I Errors with 5% and 95% Prediction Limits for DT Positive Group.**

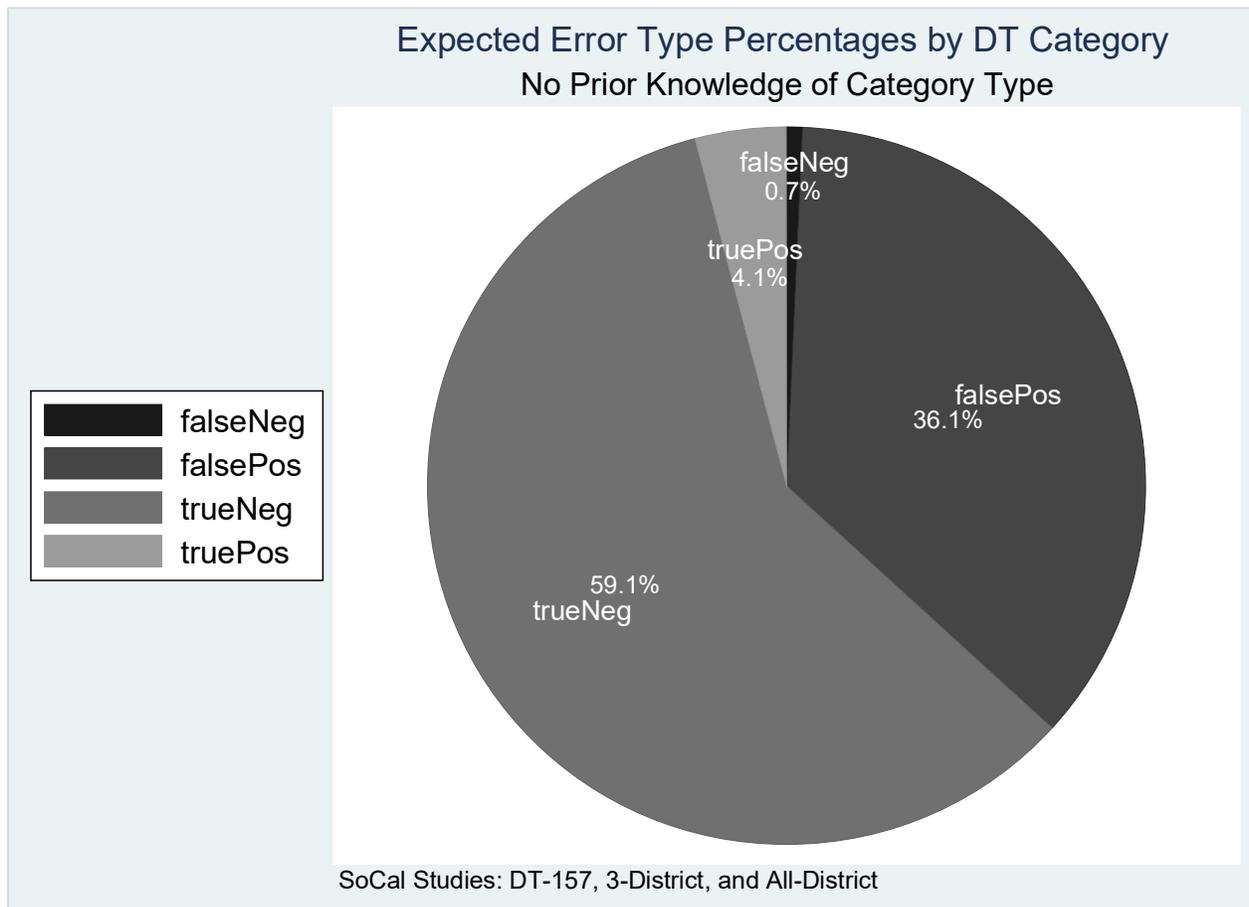
Error Type	Count	LPL%	MLV%	UPL%
False Pos	105.000	84.044	89.744	93.359
True Pos	12.000	6.641	10.256	15.956

## Independent False/True Negative Error Type II

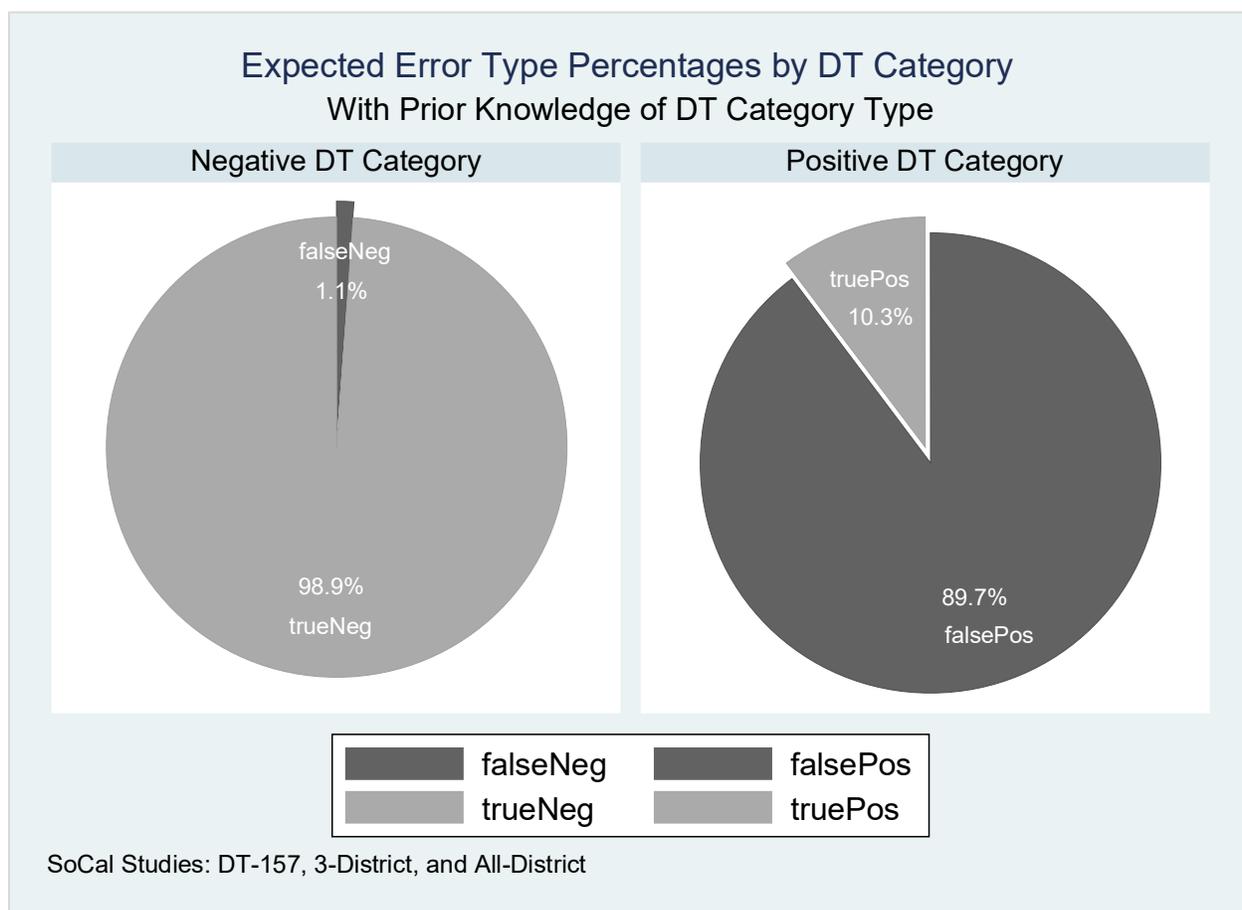
Table 26: Type II Errors with 5% and 95% Prediction Limits for DT Negative Group.

Error Type	Count	LPL%	MLV%	UPL%
False Neg	2.000	0.469	1.149	3.554
True Neg	172.000	96.446	98.851	99.531

Figure 24: Expected Decision Tree Output with No Concentration Data.



**Figure 25: Expected Decision Tree Output with Known DT Category.**



### Overall Efficiency of Decision Tree Process

The last section used Bayesian analysis to calculate the likelihoods of Type I and II errors associated with the Decision Tree criteria at establishing if a leak would be a non-hazardous large or not large leak. There was 100% leak rate testing of the 291 sample set.

Expanding this discussion outside of the sample used to characterize the DT error and output levels, another operational consideration not discussed to this point is the *efficiency* of the DT approach as to what percentage of totally screened leaks would likely result in the DT criteria being met, thereby triggering a recommendation for leak rate measurement testing.

### 2019 3-District Pilot Data Example

With the DT threshold is set at the 10 scfh level, the 2019 Pilot study had a total number of 356 leaks with surface concentration measurements. Of these, the DT was triggered for measurement 44 times. Therefore, this relates to a flow rate measurement ratio of 44 / 356 or 12.4%. In other words, when considering leak sites visited and screened with surface concentration measurements that one would expect leaks triggered by the DT process and criteria to have

approximately a 1 in 8 chance of being classified as potential non-hazardous large leak, scheduled for leak rate measurement, or prioritized for repair.

For this particular example, rather than measuring all 356 leaks to find all the large leaks; the DT process was used resulting in the requirement to measure only 1 in 8 leaks while maintaining a false negative rate of 1.1%. In summary:

- Using the DT method, 4 of the expected 7 large leaks were found by measuring the leak flow rate from 44 out of 356 leak sites. The value of 7 came from the 4 large leaks discovered from direct flow rate measurements of predicted large leaks, plus the 3 calculated from the DT false negative percent applied to the samples not predicted to be large leaks by the DT process.
- Without the DT, to find the same ratio of 4 out of the 7 large leaks, 203 leak flow rates on average would need to be measured out of the 356 leak sites.
- This means the DT efficiency increase is  $203/44 = 4.6x$  (460%) more efficient at finding the same number of large leaks when not using the DT process.
- The DT is therefore an efficient screening mechanism, with a high potential to continue to improve over the short-term full implementation period.

#### **Additional Ongoing False Negative Validation Sampling**

In order to continually confirm and refine the false negative error rate of the DT, the leak investigation process will measure an additional 59 sites (for a 90% confidence level) that the DT predicts to not be a large leak. This amounts to 1-2% of the annually encountered leaks in the field (e.g., a total of about 3,000 to 6,000 leaks per year).

## 5.10. Population Mean Leak Rate Analysis

In this and subsequent sections of this report, the emphasis will be utilizing the leak flow measurement data to create a leak emission factor based on the Decision Tree analysis process.

### Bootstrap Analysis of Field Leak Rate Data

A non-parametric bootstrap analysis [18-20] using resampling with replacement was conducted to establish the mean leak rates and a full set of mean percentiles of each of the studies as well as the combined SoCalGas studies.

The combined SoCalGas study included the overall mean leak rate and the mean leak rates for actual leak situations less than 10 scfh and greater than or equal to 10 scfh. The resample size was set to the same size as the field sample size, and the number of resamples was set to 10,000.

### Monte Carlo Analysis of Fitted Distribution (for illustrative purposes only)

Additionally, the combined SoCalGas study samples were fit to a log-normal distribution as discussed in an earlier section of this report. This was then analyzed with a Monte Carlo analysis with samples from the distribution fit set to the original field sample size to extract out mean leak rates for the same three categories as was done with the bootstrap analysis.

This limited sample size for the Monte Carlo sample increases the uncertainty in the average, since you do not leverage the central limit theorem with huge sample sizes. You have huge numbers of overall samples, but each has only the limited number of individual observations per sample. This leads to more uncertainty vs. less.

### Mean Leak Rate Analysis Results

The bootstrap mean leak rates and minimum and maximum *mean* leak rates from the bootstrap analysis are presented in Table 27.

The last three rows of the table also include the Monte Carlo analysis of the log-normal distribution fit of the sample leak rate distribution (details are in the next section and Appendix C). The bootstrap leak rates will be used as part of the emission rate calculations and are robust against non-normally distributed data.

**Table 27: Leak Rate Bootstrap Means by Study Group.**

<b>Bootstrap Leak Rate Means (10,000 Resamples)</b>			
<b>Study</b>	<b>mean(scfh)</b>	<b>min(scfh)</b>	<b>max(scfh)</b>
Natl CARB	2.484	1.098	4.756
Natl OTD	5.767	1.375	16.397
Natl WSU	1.682	0.506	4.462
SoCal DT-157	2.926	1.644	5.280
SoCal All-Dist	1.569	0.383	3.647
SoCal 3-Dist	12.043	0.826	54.358
SoCal All	4.303	1.635	12.013
SoCal LT10	1.623	1.141	2.153
SoCal GE10	57.667	14.230	194.943
SoCal L-N Fit All	4.960	2.483	26.049
SoCal L-N Fit L10	2.871	1.452	25.009
SoCal L-N Fit GE10	46.294	15.409	188.861

L-N: Log-normal fit to SoCalGas combined study is explained later in this section.  
 L10: Leak rate is less than 10 scfh  
 GE10: Leak rate is greater than or equal to 10 scfh

The 5<sup>th</sup> through 95<sup>th</sup> percentiles for the mean leak rates are presented in Table 28 for the bootstrap analysis.

**Table 28: Leak Rate (scfh) Percentiles of the Bootstrap Mean.**

<b>Pct</b>	<b>CARB</b>	<b>OTD</b>	<b>WSU</b>	<b>DT157</b>	<b>All-Dist</b>	<b>3-Dist</b>	<b>SoCal</b>	<b>SoCalL10</b>	<b>SoCal GE10</b>
5	1.783	2.838	0.803	2.247	0.909	2.124	2.331	1.421	20.786
10	1.921	3.331	0.915	2.376	1.031	2.679	2.574	1.463	24.487
15	2.020	3.721	1.050	2.463	1.114	4.824	2.790	1.493	30.401
20	2.100	4.011	1.150	2.535	1.187	5.347	3.015	1.517	33.143
25	2.166	4.284	1.225	2.604	1.249	6.058	3.246	1.537	39.146
30	2.229	4.551	1.303	2.667	1.309	8.253	3.433	1.555	43.303
35	2.294	4.808	1.377	2.729	1.365	8.625	3.606	1.573	45.360
40	2.350	5.073	1.457	2.785	1.416	8.986	3.781	1.590	47.597
45	2.410	5.325	1.539	2.838	1.475	9.588	3.945	1.605	51.525
50	2.460	5.570	1.613	2.894	1.529	11.469	4.121	1.622	55.236
55	2.512	5.810	1.694	2.951	1.587	11.876	4.289	1.637	57.394
60	2.572	6.058	1.769	3.007	1.644	12.363	4.488	1.652	60.164
65	2.632	6.369	1.858	3.068	1.707	14.587	4.697	1.669	66.590
70	2.697	6.689	1.952	3.136	1.779	15.173	4.924	1.687	69.620
75	2.776	7.026	2.058	3.207	1.849	15.766	5.167	1.706	72.679
80	2.856	7.411	2.186	3.288	1.931	18.171	5.452	1.727	79.532
85	2.963	7.841	2.327	3.394	2.026	19.000	5.793	1.752	83.776
90	3.072	8.444	2.499	3.523	2.152	21.846	6.197	1.785	93.578
95	3.259	9.397	2.786	3.719	2.358	25.268	6.894	1.834	105.665

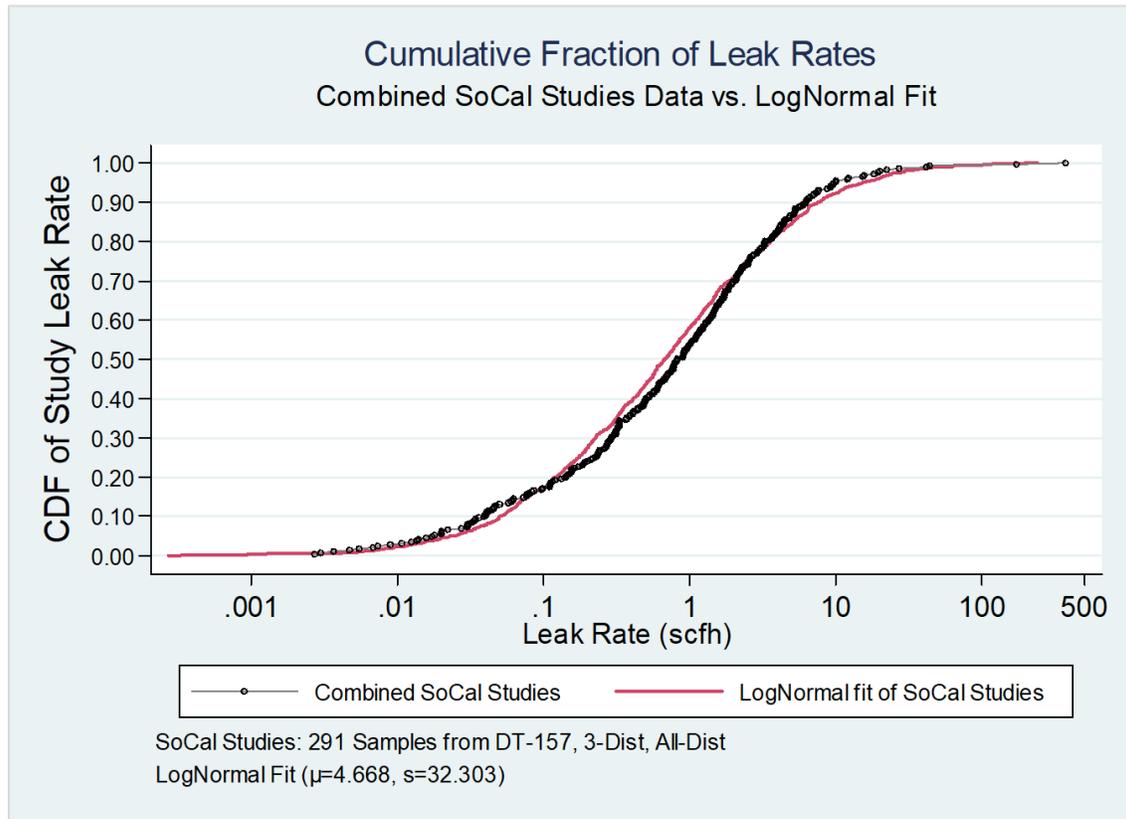
## Log-Normal Distribution Fit and Monte Carlo Analysis

The log-normal distribution fit of the SoCalGas combined study leak rates is shown in Figure 26.

The black points are the 291 empirical data points from the field. The red line is the fit of a log-normal distribution with the two parameters noted on the plot.

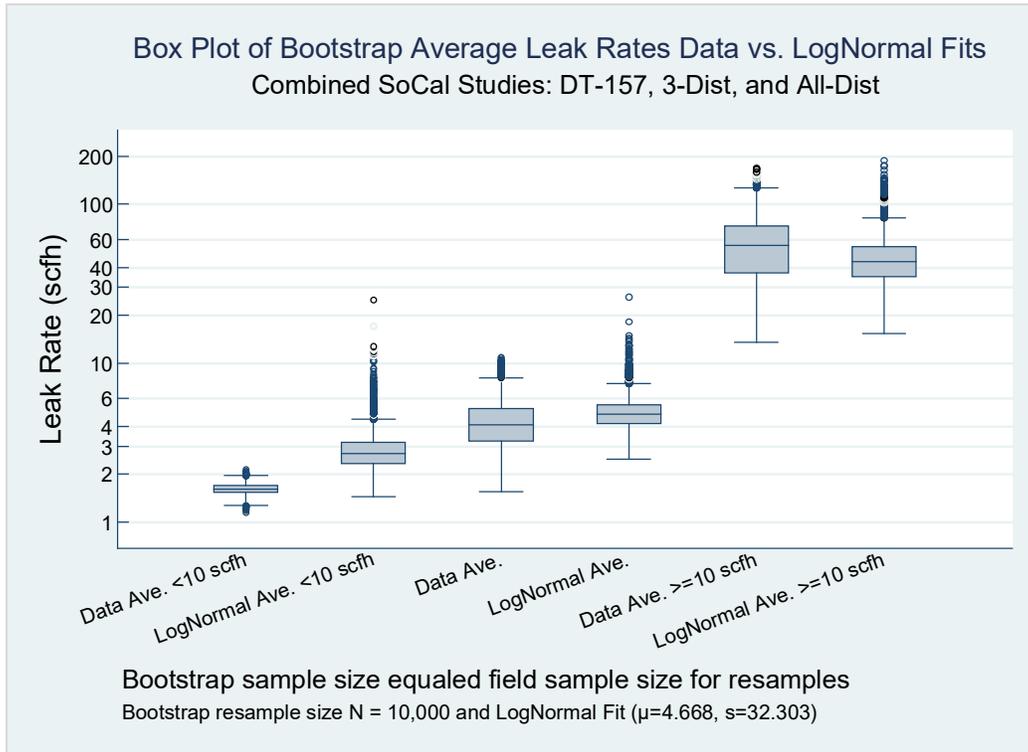
The fit is very strong with three relatively tight intersections between 0.01 and 100 scfh. A summary of the log-normal distribution is in Appendix C.

**Figure 26: Leak Rate Cumulative Fractions of Combined SoCalGas Studies and L-N Fit.**



A summary box plot of the bootstrap 10,000 resamples (at the original sample sizes) from the field leak rate data and the separate Monte Carlo simulations (also sampled at the original field sample size) from the fitted distribution is shown in Figure 27.

**Figure 27: Bootstrap Mean Leak Rate Box Plots of Field Data and Log-Normal Fit.**



One can see that the log-normal data is conservative in the low and mid-range, and about the same as the bootstrapped field data in the high range.

The log-normal (or other appropriate fit such as log-gamma) fit with a Monte Carlo analysis can therefore provide an alternative to the bootstrap resampling of the actual field leak rate data until a large enough sample size is available from the field to run the bootstrap analysis.

Once a significant bootstrap sample size is established one could then shift to the results of that analysis vs. fitting the distribution and running a Monte Carlo analysis.

## 6. Emission Factor Development and Application

---

### 6.1 Development of a Company Specific Emission Factor

As stated earlier, the objective of this study was to develop a method for flagging large leaks for cost-effective measurement and repair to minimize system-wide methane leakage rates. So, if a company can reduce its number of higher emitting non-hazardous leaks, it can reduce actual emissions *and* more accurately estimate the reduction.

With the information assembled in the previous sections, a SoCalGas-specific non-hazardous leak emission factor can be reliably developed based on sound statistical and probabilistic sampling and analysis.

#### Input Information

To construct an accurate emission factor, the following steps were taken (three decimal places are listed to prevent round-off errors in calculations):

1. From Table 27, use three specific bootstraps (log-normal fits if there is not enough samples to execute a bootstrap analysis, e.g. 30 random samples) of the mean leak rates from the SoCalGas studies, for the:
  1. Entire sample set (**ALL**): 4.303 scfh
  2. Samples < 10 scfh (**L10**): 1.623 scfh
  3. Samples ≥ 10 scfh (**GE10**): 57.667 scfh
2. Use the Bayesian Error Table most likely value (MLV) proportions for true and false negatives and positives from Table 25 and Table 26:
  1. False Positives (**FP**): 89.744%
  2. True Positives (**TP**): 10.256%
  3. False Negatives (**FN**): 1.149%
  4. True Negatives (**TN**): 98.851%

## Calculation of Distinct Emission Factors

The emission factors are calculated by properly combining the information above for two situations, when the Decision Tree criteria is met (predicting greater than or equal to 10 scfh leak rates) and when it is not met (predicting less than 10 scfh leak rates) as follows:

1. For DT met (**≥ 10 scfh** prediction):
  - = [True Positive MLV% (**TP**)] x [Mean Leak Rate for ≥ 10 scfh (**GE10**)] + [False Positive MLV% (**FP**)] x [Mean Leak Rate for < 10 scfh (**L10**)]
  - = 10.256/100 x 57.667 scfh + 89.744/100 x 1.623 scfh
  - = **7.37 scfh**
2. For DT not met (**< 10 scfh** prediction):
  - = [True Negative MLV% (**TN**)] x [Mean Leak Rate for < 10 scfh (**L10**)] + [False Negative MLV% (**FN**)] x [Mean Leak Rate for ≥ 10 scfh (**GE10**)]
  - = 98.851 /100 x 1.623 scfh + 1.149/100 x 57.667 scfh
  - = **2.27 scfh**

Note that the bootstrap mean of the leak rate for samples greater than or equal to 10 scfh was used conservatively in the above calculation. One can shift from this value to a bootstrap mean of the leak rate of the False Negatives of the sample data when enough of these values are obtained. There are currently only two occurrences of False Negatives, so a population mean cannot be obtained.

3. If no concentration measurements were taken, i.e. no application of Decision Tree or leak rate measurements, then one should use the entire sample set bootstrap mean:
  - = **4.30 scfh**
4. If one has an actual leak rate measurement, then use that measurement.

## 6.2. Table of Emission Factors

The SoCalGas-specific emission factors calculated above are summarized in Table 29 below.

**Table 29: Table of SoCalGas Company Specific Emission Factors by DT Grouping.**

EF Category	EF (scfh)
Combined All Case Ave EF	4.30
DT Not Triggered Ave EF	2.27
DT Triggered Ave EF	7.37

It is worth noting that these DT related emission factors are conservative due to the nature of properly accounting for false negatives. For example, if one were to take the straight average of the All District Study from the 78 samples, the single emission factor would be 1.58 scfh.

### 6.3. Carrying Uncertainty Through to the Emission Factor Calculations

Additional steps are necessary to properly carry through the uncertainty related to the average (i.e., expected or baseline) emission factor and provide confidence limits at a selected confidence level for the EF's.

To do this, one would run Monte Carlo analysis by drawing from the bootstrap average leak rate population distributions of the appropriate data set and category of leak rate (large and not large) and then weight those by the Bayesian proportions for those categories. This would be done thousands of times, picking the average leak flow rates and the associated Bayesian proportions from those distributions and calculating (thousands of times) the associated emission factors.

This would provide a full distribution of the emission factors for each category and then one could select the confidence level of choice (e.g., 95%) to generate the confidence interval around the average emission factors.

However, one would still use the expected (average) value of the emission factors in practice, but the confidence bands would help establish the level of uncertainty in those values.

This will be the next step once SoCalGas collects additional samples from the ongoing implementation of this approach. As was shown in an earlier section, it is desired to get a statistically significant sample for false negatives, so that data can be used for the associated average leak flow rate vs. the much more conservative measure currently being used, which is the average of actual leak flow rates above the 10 scfh threshold point.

## 6.4. Scenarios of EF Application

The scenarios that could be encountered in the field for leak repair are listed below with the guidance on what emission factor to use and when to use them.

**Table 30: Table of Emission Factors to use for Field Situations.**

<b>Situation Number</b>	<b>Field Situation Description</b>	<b>Emission Factor</b>
1	Measured concentration triggers DT < 10 scfh category & leak rate is not measured (which would be the typical situation) - Use DT Not Triggered Ave EF	<b>2.27 scfh</b>
2	Measured concentration DT ≥ 10 category & leak rate is not measured (used when leak rate cannot be measured, such as leaks quickly repaired or when leak is in a remote location) - Use DT Triggered Ave EF	<b>7.37 scfh</b>
3	Leak repaired and no concentration or leak rate measurements - Use Combined All Case Ave EF	<b>4.30 scfh</b>
4	Measured concentration(s) trigger DT >10 category & then leak rate measured and actual leak rate is < 10 scfh - Use the actual leak rate measurement for the emission factor	<b>Use actual leak rate measurement</b>
5	Measured concentration(s) trigger DT >10 category & then measure and actual leak rate is ≥ 10 scfh - Use the actual leak rate measurement for the emission factor	<b>Use actual leak rate measurement</b>

## 7. Summary of Results and Conclusions

---

### Summary of Results

The national studies compared well with the SoCalGas studies. The upper and lower 95% percentiles for leak rate and the median and means of these two groups are similar.

Two of five SoCalGas sample sets were known to contain sample bias as well as being an order of magnitude in size smaller than the other three. These were analyzed in this report to show how bias might appear during analysis, and they were not included in the ultimate combined data set.

The non-hazardous leak rate values from the SoCalGas combined data set was analyzed for unexplainable outliers or extreme values and was log transformed, resulting in a normally distributed data set. Upon review of the extreme values, all of them were deemed as sound data points and not errors or anomalous values. The log-normal transformation of the leak rate data permitted a variety of statistical regression tools to be appropriately leveraged.

A series of regression and probabilistic analysis were conducted on the data set. A key finding was that when the samples sizes would support categorical analysis that there was no significant sensitivity of the leak rate means to geographic districts of the leak or the year that the leak was detected.

An analysis of the field methane concentration vs. measured leak rates was done by Decision Tree methane concentration threshold category. The regression analysis of the mean leak flow rate vs. methane concentration showed the expected upward trend for the average values. The concentration threshold intersection with the established 10 scfh “Large” vs. “Not Large” flow rate threshold was within the 95% confidence interval of the regression model or above and to the left (a conservative situation) of the predictive margin plots.

A Bayesian probabilistic analysis was conducted of the Decision Tree threshold performance. This resulted in a true/false positive/negative error table. The Decision Tree thresholds correctly assigned not large leak situations 98.9% of the time, i.e. true negatives with a 95% prediction interval of 98.9% to 99.5%. Likewise, the Decision Tree had a false negative (Type II error) of only 1.1% with a 95% prediction interval of 0.47% to 3.6%.

The leak rate data was bootstrapped 10,000 times with replacement with a re-sample size equal to the field data sample size. This analysis provided the overall mean leak rate, as well as the mean leak rates for less than 10 scfh leakers and greater than or equal to 10 scfh leakers - all from the empirical data. The bootstrap analysis provided a full set of percentiles for the actual mean leak rates which allows one to establish confidence intervals for the mean values at any desired confidence level.

The leak rate data was fit to a log-normal distribution as well, and this fit was used to conduct a Monte Carlo analysis of the mean leak rates as was performed with the bootstrap analysis using

the actual field leak rate data. The same re-sample and sample sizes were used as was done with the bootstrap analysis to properly propagate the uncertainty through the analysis. The result showed the two approaches were very similar, with the Monte Carlo of the log-normal distribution fit being conservative in the low- to mid- leak rate ranges and about the same in the high-leak rates.

A set of emission factors based on the Decision Tree categorization were calculated by combining the mean leak rates with their corresponding expected percentiles (in a weighted manner) from the Decision Tree error table. It was noted that the Decision Tree derived emission factors were conservative (higher) than one would have obtained from a straight average of the empirical data from the All District Study of the SoCalGas system. This is due to the Bayesian analysis properly accounting for false negatives in the Decision Tree process.

A calculation of the efficiency of the process was done using the 2019 3-District Pilot study which had a total number of 356 screened leaks with surface concentration measurements. Of these, the DT was triggered for measurement 44 times. Therefore, this relates to a flow rate measurement ratio of 44 / 356 or 12.4%. In other words, when considering leak sites visited and screened with surface concentration measurements that one would expect leaks triggered by the DT process and criteria to have approximately a 1 in 8 chance being classified as potential non-hazardous large leak, scheduled for leak rate measurement, or prioritized for repair.

For this particular example, rather than measuring all 356 leaks to find all the large leaks; we used the DT process resulting in the requirement to measure only 1 in 8 leaks while maintaining a false negative rate of 1.1%. In summary:

- Using the DT method, 4 of the expected 7 large leaks were found by measuring the leak flow rate from 44 out of 356 leak sites.
- Without the DT, to find the same ratio of 4 out of the 7 large leaks, 203 leak flow rates on average would need to be measured out of the 356 leak sites.
- This means the DT efficiency increase is  $203/44 = 4.6x$  (460%) more efficient at finding the same number of large leaks when not using the DT process.
- The DT is therefore an efficient screening mechanism, with a high potential to continue to improve over the short-term full implementation period.

## **Conclusions**

SoCalGas conducted a statistically sound study of underground pipeline leaks using random samples as well as well-proven field leak concentration and flow rate measurement techniques to calculate SoCalGas company-specific natural gas emission factors for buried distribution system non-hazardous leaks.

The developed Decision Tree approach of using concentration measurements with thresholds to establish large and not large non-hazardous leaks was successful as shown by a 98.9% true negative value associated with predicted leak and actual leak rates.

The inferred population mean leak rates were combined with the associated Decision Tree performance percentages to calculate appropriately weighted emission factors for large and not large non-hazardous leaks.

This allows the assignment of emission factors for the not large non-hazardous leaks that would not have leak rate flow measurements performed on them as well as any Decision Tree classified large non-hazardous leaks that did not have leak rate flow measurements performed.

The approach will be further refined and improved by continuing to:

- Collect field data leading to lower uncertainty, i.e. tighter confidence intervals around leak and Decision Tree performance metrics;
- Perform random checks for false negatives to identify possible upset conditions in expected leak rates, e.g. from a change in system performance and/or environmental stressors; and
- Analyze and adjust the Decision Tree thresholds or even add new thresholds to further increase the method's predictive accuracy and/or increase process efficiency to continuously improve the cost-effectiveness of the approach, overall process for detection, and repair of large flow system leaks to minimize natural gas emissions.

## Appendix A: Surface Measurements of Underground Leak Flow Rate

The approach employed a dynamic surface enclosure where ambient air is drawn through a well-mixed chamber at a constant measured rate. Methane emitted from the surface and mixed into the chamber air is sampled in the exhaust. The methane emission rate is calculated from the measured air flow rate through the chamber, measured inlet, and exhaust methane concentrations. This surface measurement technique for underground leaks have been validated and used extensively by many research teams, including GTI, in past research efforts [1, 4, 30].

The equipment required for this method includes an enclosure, a high-volume sampler (such as Bacharach's Hi Flow Sampler), and a combustible gas indicator (CGI), see Figure 28.

The Hi Flow sampler is a portable, battery-powered instrument designed to quantify methane emission rates from leaking components common to natural gas operations. When using the Hi Flow sampler, a robust validation procedure was followed to eliminate measurement issues as suggested by prior studies [31, 32]. In addition, the methane measurement instrument used in conjunction with the high-flow sampler was calibrated according to manufacturer requirements.

The use of the surface measurement method is suitable when leak locations are known and are accessible on foot. The underground leak first has to be identified using a screening instrument such as a handheld leak survey instrument (e.g. the DP-IR) that can be used to map out the area on the surface with elevated methane concentrations.

**Figure 28: Quantifying surface flux rate of an underground emission.**

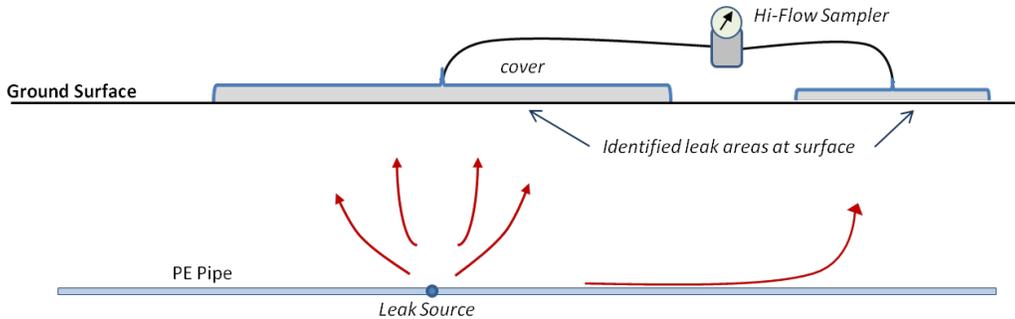


Using the "enclosure/chamber method. The high-flow device is housed in the backpack.

Once the leak area is demarcated, if it is larger than the footprint of the enclosure, then it is segmented into a grid and each square segment is then measured with the enclosure to capture the surface expression of the area. A picture and schematic of this measurement method is shown in Figure 29 below. The total leak rate is the sum of the individual grid measurements.

During measurement, the top of the enclosure is attached to a high-volume sampler that pulls in air from the enclosed volume at a high flow rate. To calculate the leak flux rate under the enclosure, the sampling rate is multiplied with methane concentration of the sampled air; this is measured by a built-in methane sensor or by a separate combustible gas indicator (CGI). The built-in methane sensor has an accuracy of 0.02% methane which gives the device a sensitivity to detect natural gas at a leak rate of 0.6 scfh [33]. The unit also corrects temperature compensates automatically to 60 F. In order to improve the sensitivity, a CGI with low parts-per-million (ppm) methane sensitivity is placed at the outlet of the high-flow device.

**Figure 29: Schematic of surface chamber measurements with the Hi-Flow sampler.**



# Appendix B: Statistical and Probabilistic Analysis Details and Supplemental Analysis

This appendix contains supplemental statistical and probabilistic analysis details that are summarized in the body of the report.

## Linear Regression Residual Analysis and Regression Diagnostics

The following diagnostics were analyzed for the SoCalGas study regressions to confirm regression assumptions and to look for influential, outlier, or extreme values requiring further review and/or explanation or exclusion.

The residuals were from the full linear regression where the dependent variable is the log(10) of the leak flow rate (scfh), and the independent, categorical variables are the three SoCalGas studies.

The plots of the analysis are customized. Instead of just listing table data of these diagnostic measures, the values were scatter plotted, and the ID labels (observations) were paired to the diagnostic results *after* the regression analysis and plotted with the ID as the x-axis in some of the cases. This was done to allow one to quickly identify observations that should be focused on for further review and spot trends. The regression was not tied to the observation (sample ID) in any way.

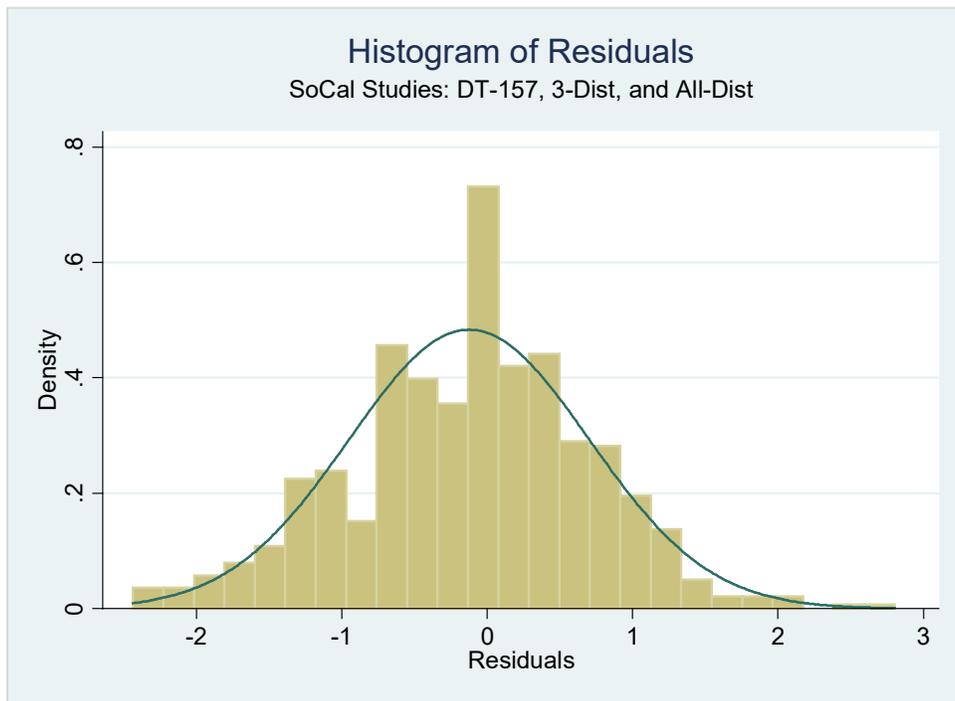
1. **Exogeneity:** Assumes that given an independent set of variables, one can account for any error in the linear regression model. Stated a different way, one must have a sound causal model with factors for each casual influence on the dependent variable. There is no statistical test, and it was therefore designed into the regression model through cause and effect analysis/modeling (causal analysis).
2. **Random Sampling:** This was confirmed and discussed prior to sampling or regression. Two sample-biased studies were not included primarily due to confirmation bias and for small sample sizes. There is no statistical test for this, and it was therefore designed into the study.
3. **Linearity in Parameters:** Since the independent variables within this study are not metric, linearity in parameters were not compared within this report.
4. **Multicollinearity:** When a large number of independent (control) variables exists, there should not be strong correlations between them, which can inflate standard errors of estimated coefficients. This is not an issue with this analysis, since the datasets are from studies that are mutually independent of each other.
5. **Heteroscedasticity and Normal Distribution of Residuals:** Assumes that variance of the residuals is constant - if it is not, then this is known as heteroscedasticity. In this study, the residuals between the predicted and actual values were analyzed, demonstrated a symmetric

and normal distribution across zero (as plotted in Figure 30). Likewise, Figure 31 is a scatter plot comparing the residuals vs. the study (independent variable) with fitted values plotted as a line. In addition, a locally weighted (Lowess) line fit was conducted and appended onto Figure 31. The two lines are flat, nearly parallel to each other, and centered around zero. Note the values above and below these lines are nearly equally spread in numbers and density.

6. **Influential Observations - DFBETA:** This diagnostic is used to measure the influence of a single observation per each metric or categorical variable. A rule of thumb is that DFBETA might be a problem when an observation's absolute value of DFBETA is  $> 2/\sqrt{N}$ , where N is the total number of samples (observations). In this case, the value is 0.12. As can be seen from Figure 32, nearly all the cases are below this threshold number. To analyze values above the threshold number, Cook's distance and Leverage measures of influence were also utilized.
7. **Influential Observations - Cook's Distance:** This measure detects strange patterns and unusual variable combinations. Upon looking at DFBETA and Cook's distance, observation ID 97 and ID 121 are more influential than other observations (Figure 33). This makes sense, since these two observations are the two highest leak rates at 373 and 172 scfh (in round numbers). However, these are not problematic or erroneous outliers and *should* be considered influential, since they are in fact within the far-right tail of the expected leak rate distribution of the sample.
8. **Influential Observations - Leverage:** Finally, a dual plot of leverage vs. residuals is a very useful plot to analyze. If the *residual* of a case is high, this means that the regression would calculate a result that is quite off from the real outcome. Therefore, the residual is related to the *dependent* variable (leak rate in this case). A high *leverage* of a case means that the constellations of independent variables of a certain case are so extreme or uncommon that they influence the final result over proportionally. Therefore, the leverage is related to the independent variables of a case. The two red lines in the graph show the means for both residuals and leverages. Looking at the dual plot in Figure 34, one can see two observations in the upper right quadrant which are again the two highest leak rate samples as noted earlier. There are also three observations in the lower right quadrant which happen to be the three lowest leak rate readings. Since these observations are correct, it was decided to retain them in the analysis.

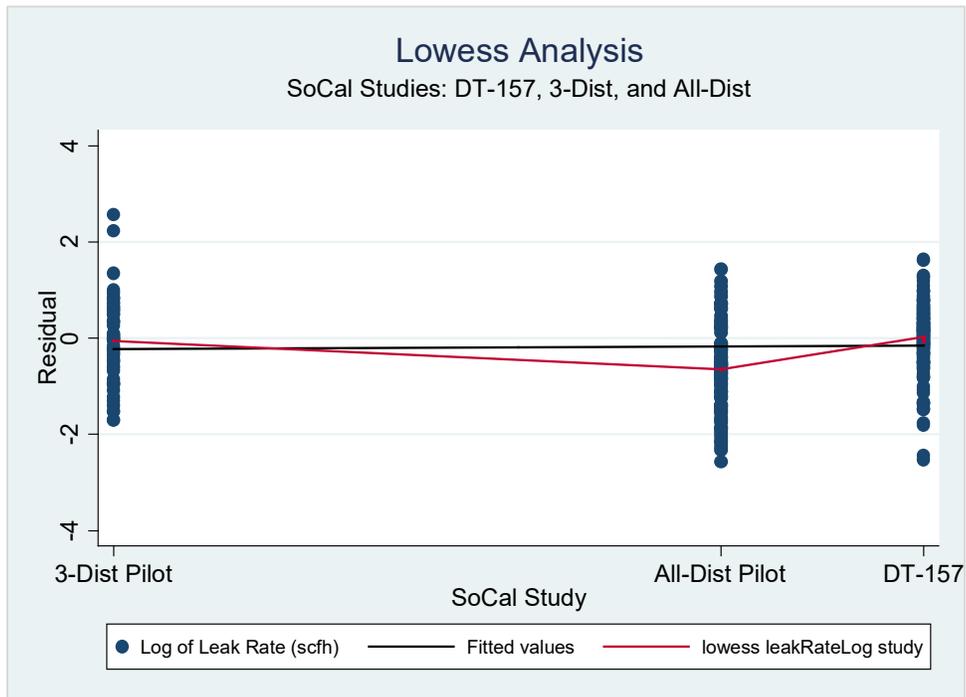
## Residual Distribution

Figure 30: Histogram Diagnostic Plot of LR Leak Rate Residuals for Three SoCalGas Studies.



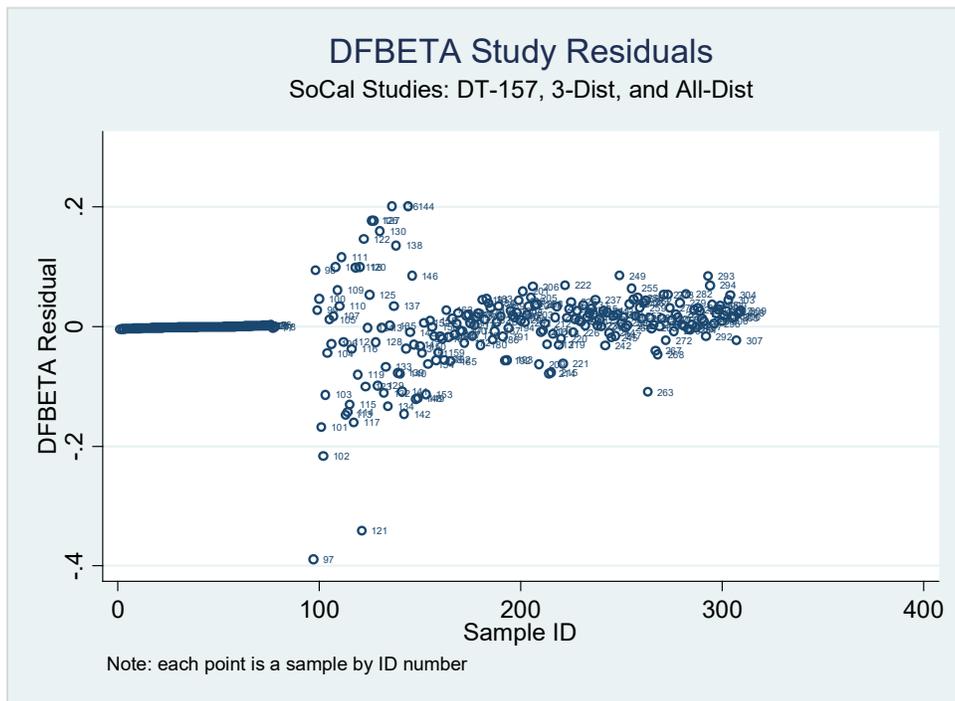
## Lowess

Figure 31: Lowess LR Diagnostic of Leak Rate Residuals by SoCalGas Study.



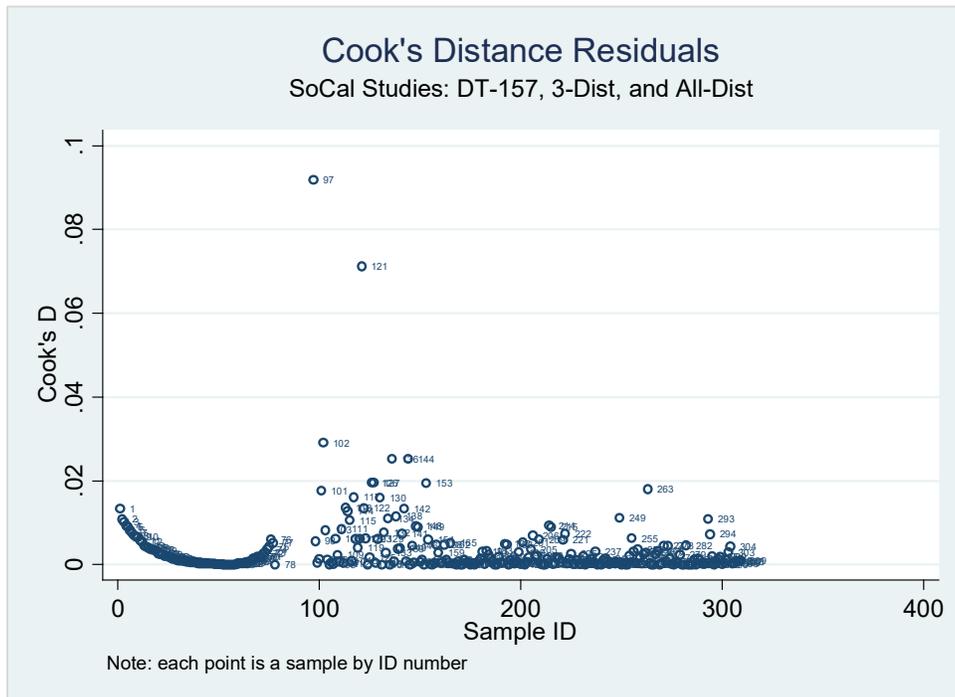
## DFBETA

Figure 32: DFBETA LR Diagnostic of Leak Rate for Three SoCalGas Studies.



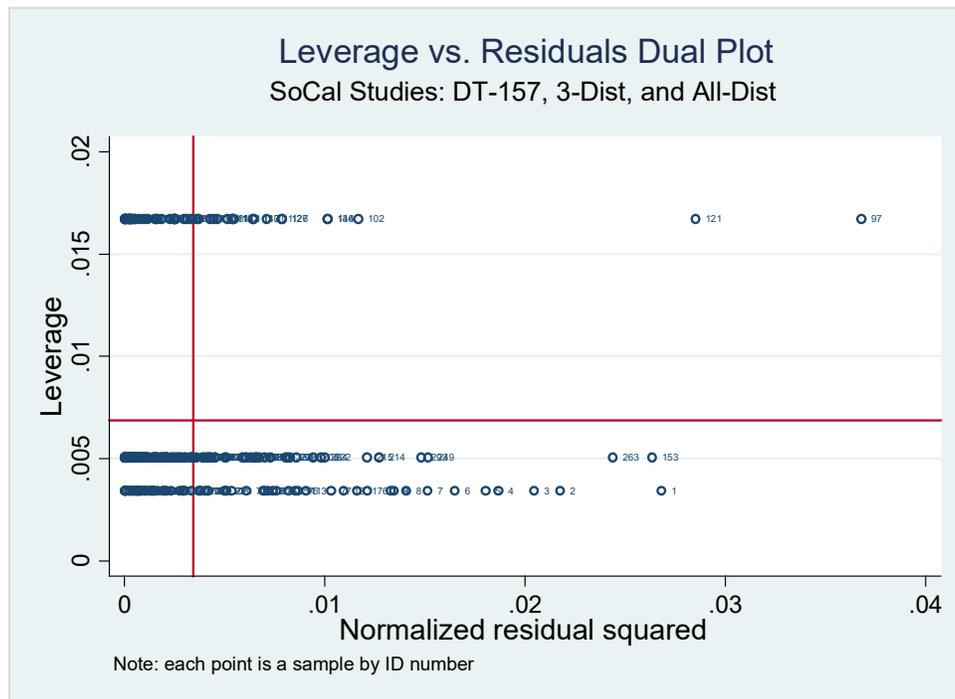
## Cook's Distance

Figure 33: Cook's Distance LR Diagnostic of Leak Rate for Three SoCalGas Studies.



## Leverage

Figure 34: Leverage LR Diagnostic of Leak Rate for Three SoCalGas Studies.



## Bayesian Monte Carlo Markov Chain (MCMC) Models of Sample Leak Rates

A non-parametric Bayesian Monte Carlo Markov Chain (MCMC) [25] regression analysis was conducted. Two variations of random *sampling* were used, the random walk Metropolis-Hasting (MHS) [26, 27] method as well as the more robust Gibbs (GS) [28] method. In both cases, 35,000 iterations were used with a 5,000 iteration burn-in run, resulting in an incorporated Monte Carlo sample size of 30,000. In both cases, the prior distribution for the log(10) leak rate distribution was "uniformed", i.e. a flat/uniform prior. The sigma prior was assumed as a conservative gamma function. The same set up of dependent and independent variables was used as with the traditional linear regression and ANOVA analysis discussed in this appendix and the body of the report.

### Metropolis-Hastings Sampling (MHS)

The results of the MCMC-MHS are shown in Table 31 below. The form is very similar in output to the standard regression outputs already explained. However, this is for convenience of comparison, since the methods are completely different, and this analysis uses Bayesian linear regression. The MHS method is about 33% efficient. Therefore, approximately 10,000 of the samples were utilized with the others being rejected.

From the table, the MCMC-MHS produces very similar coefficients and associated 95% credible interval as compared to the LR coefficients (means) and 95% confidence interval. This does not come as a surprise, since the regression assumptions were met, and the dependent variable (log(10) of leak rate) was normally distributed.

**Table 31: Bayesian MCMC(MHS) of Leak Rate Means - Three SoCalGas Studies.**

Model summary						
-----						
Likelihood:						
leakRateLogSoCal ~ regress xb_leakRateLogSoCal,sigma2)						
Priors:						
leakRateLo~l:i.study _cons ~ 1 (flat)						(1)
sigma2 ~ igamma(.01,.01)						
-----						
(1) Parameters are elements of the linear form xb_leakRateLogSoCal.						
Bayesian linear regression			MCMC iterations =	35,000		
Random-walk Metropolis-Hastings sampling			Burn-in =	5,000		
			MCMC sample size =	30,000		
			Number of obs =	291		
			Acceptance rate =	.3335		
			Efficiency: min =	.06429		
			avg =	.1081		
			max =	.2204		
Log marginal-likelihood = -360.84393						
-----						
	Mean	Std. Dev.	MCSE	Median	Equal-tailed [95% Cred. Interval]	
-----						
leakRateLogSoCal						
study						
3DisPilot	(base)					
3DisPilotLowSpec	(omitted)					
AllDisLIRP	(omitted)					
AllDisPilot	-.5893627	.1408558	.002931	-.5869901	-.8674918	-.3101206
DT157Pilot	.0853107	.1241679	.002827	.0824296	-.1555874	.3314519
Natl_CARB_GTI	(omitted)					
Natl_OTD_GTI	(omitted)					
Natl_WSU_EDF	(omitted)					
_cons	-.063096	.1081477	.002349	-.0618196	-.2803022	.1406578
-----						
sigma2	.659823	.0556692	.000685	.657526	.5587767	.7786351

## Gibbs Sampling (GS)

A very similar analysis was repeated with the only difference being that Gibbs sampling was used instead of Metropolis-Hastings Sampling. Gibbs sampling is even more efficient (note a 99% efficiency), but it does take more computing power and time to complete. The results are in Table 32 below.

**Table 32: Bayesian MCMC(GS) of Leak Rate Means for Three SoCalGas Studies.**

Model summary						
-----						
Likelihood:						
leakRateLogSoCal ~ normal(xb_leakRateLogSoCal,sigma2)						
Priors:						
leakRateLo~l:i.study _cons ~ normal(0,10000)						(1)
sigma2 ~ igamma(.01,.01)						
-----						
(1) Parameters are elements of the linear form xb_leakRateLogSoCal.						
Bayesian linear regression			MCMC iterations =	35,000		
Gibbs sampling			Burn-in =	5,000		
			MCMC sample size =	30,000		
			Number of obs =	291		
			Acceptance rate =	1		
			Efficiency: min =	.9682		
			avg =	.9854		
			max =	1		
Log marginal-likelihood = -377.39597						
-----						
	Mean	Std. Dev.	MCSE	Median	Equal-tailed [95% Cred. Interval]	
-----						
leakRateLogSoCal						
study						
3DisPilot	(base)					
3DisPilotLowSpec	(omitted)					
AllDisLIRP	(omitted)					
AllDisPilot	-.5896894	.1425934	.000834	-.5907453	-.8732188	-.3109355
DT157Pilot	.083985	.1267482	.000732	.0830794	-.1636278	.3327636
Natl_CARB_GTI	(omitted)					
Natl_OTD_GTI	(omitted)					
Natl_WSU_EDF	(omitted)					
_cons	-.0621054	.1084323	.000626	-.0618162	-.2748983	.1485988
-----						
sigma2	.661009	.0557936	.000327	.6578556	.5605393	.7781685

Diagnostics for one of the coefficients were selected for inclusion in the report, see Figure 35 for the MCMC(MHS) model and Figure 36 for the MCMC(GS) model. In both cases, the traces are flat and well spread out, the histogram and densities are symmetric and consistent, and autocorrelation decreases quickly or is non-existent; these are all excellent attributes of the residuals.

Figure 35: Diagnostics for Bayesian MCMC(MHS) of Leak Rate Means - SoCalGas Studies.

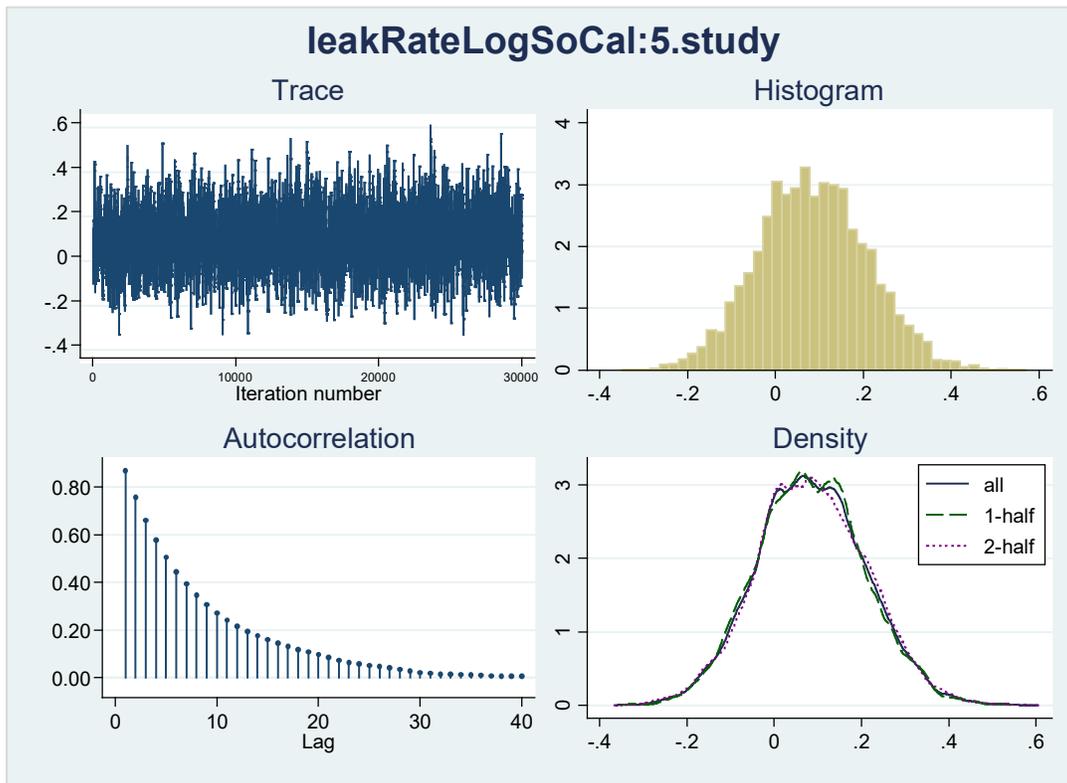
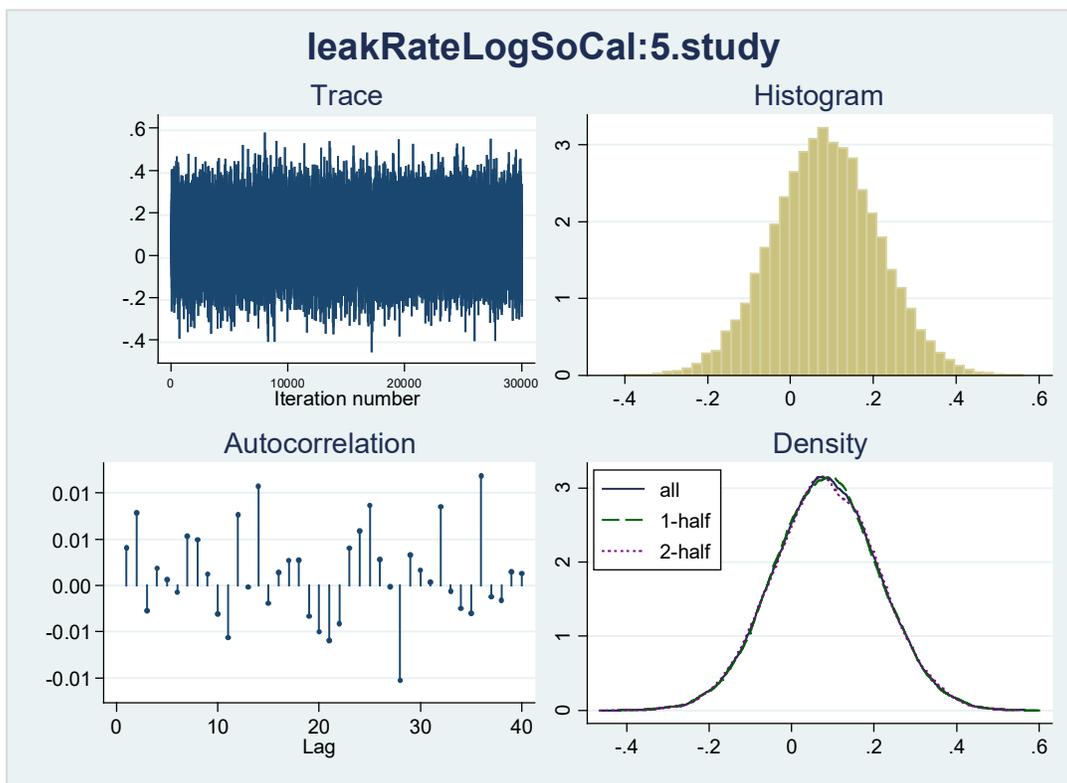


Figure 36: Diagnostics for Bayesian MCMC(GS) of Leak Rate Means - SoCalGas Studies.



## ANOVA and Pairwise comparison of National and SoCalGas Studies

The ANOVA results of a log(10) means analysis by individual study is shown in Table 33 below. The analysis was done across national industry and SoCalGas studies.

**Table 33: ANOVA of Individual National and SoCalGas Study Leak Rate Means.**

Summary of Log of Leak Rate (scfh)			
study	Mean	Std. Dev.	Freq.
3DisPilot	-.06279222	.89898228	56
3DisPilotLS	-.65094811	.55460313	8
AllDisLIR	.07601364	.65975028	10
AllDisPil	-.65257563	.92139692	78
DT157Pilo	.02154301	.71202596	157
Natl_CARB	-.01830655	.64701923	76
Natl_OTD_	.16791675	.61032353	62
Natl_WSU_	-.61104802	.75895091	212
Total	-.26707906	.81914558	659

Analysis of Variance					
Source	SS	df	MS	F	Prob > F
Between groups	70.8808391	7	10.1258342	17.79	0.0000
Within groups	370.636818	651	.56933459		
Total	441.517657	658	.670999479		

Bartlett's test for equal variances:  $\chi^2(7) = 20.9465$  Prob> $\chi^2 = 0.004$

The pairwise log(10) leak rate means comparison by individual study is shown in Table 34 below.

**Table 34: PW Comparison by National/SoCalGas Study Pair for Leak Rate Means.**

study	Contrast	Std. Err.	Unadjusted		[95% Conf. Interval]	
			t	P> t		
3DisPilotLowSpec vs 3DisPilot	-.5881559	.2851903	-2.06	0.040	-1.14816	-.028152
AllDisLIRP vs 3DisPilot	.1388059	.259037	0.54	0.592	-.3698429	.6474547
AllDisPilot vs 3DisPilot	-.5897834	.1321584	-4.46	0.000	-.8492916	-.3302752
DT157Pilot vs 3DisPilot	.0843352	.1174437	0.72	0.473	-.146279	.3149494
Natl_CARB_GTI vs 3DisPilot	.0444857	.1328832	0.33	0.738	-.2164457	.305417
Natl_OTD_GTI vs 3DisPilot	.230709	.1391025	1.66	0.098	-.0424347	.5038527
Natl_WSU_EDF vs 3DisPilot	-.5482558	.1133677	-4.84	0.000	-.7708662	-.3256454
AllDisLIRP vs 3DisPilotLowSpec	.7269618	.357911	2.03	0.043	.0241625	1.429761
AllDisPilot vs 3DisPilotLowSpec	-.0016275	.2801178	-0.01	0.995	-.551671	.5484159
DT157Pilot vs 3DisPilotLowSpec	.6724911	.2734834	2.46	0.014	.1354751	1.209507
Natl_CARB_GTI vs 3DisPilotLowSpec	.6326416	.2804605	2.26	0.024	.0819253	1.183358
Natl_OTD_GTI vs 3DisPilotLowSpec	.8188649	.2834601	2.89	0.004	.2622584	1.375471
Natl_WSU_EDF vs 3DisPilotLowSpec	.0399001	.2717579	0.15	0.883	-.4937277	.5735279
AllDisPilot vs AllDisLIRP	-.7285893	.2534416	-2.87	0.004	-1.226251	-.2309277
DT157Pilot vs AllDisLIRP	-.0544706	.246089	-0.22	0.825	-.5376946	.4287533
Natl_CARB_GTI vs AllDisLIRP	-.0943202	.2538202	-0.37	0.710	-.5927253	.4040849
Natl_OTD_GTI vs AllDisLIRP	.0919031	.2571309	0.36	0.721	-.4130028	.596809
Natl_WSU_EDF vs AllDisLIRP	-.6870617	.24417	-2.81	0.005	-1.166518	-.2076058
DT157Pilot vs AllDisPilot	.6741186	.1045251	6.45	0.000	.4688716	.8793657
Natl_CARB_GTI vs AllDisPilot	.6342691	.1216158	5.22	0.000	.3954625	.8730757
Natl_OTD_GTI vs AllDisPilot	.8204924	.1283822	6.39	0.000	.5683993	1.072585
Natl_WSU_EDF vs AllDisPilot	.0415276	.0999235	0.42	0.678	-.1546836	.2377388
Natl_CARB_GTI vs DT157Pilot	-.0398496	.1054399	-0.38	0.706	-.246893	.1671938
Natl_OTD_GTI vs DT157Pilot	.1463737	.1131775	1.29	0.196	-.0758633	.3686108
Natl_WSU_EDF vs DT157Pilot	-.632591	.0794473	-7.96	0.000	-.7885949	-.4765871
Natl_OTD_GTI vs Natl_CARB_GTI	.1862233	.1291281	1.44	0.150	-.0673345	.4397811
Natl_WSU_EDF vs Natl_CARB_GTI	-.5927415	.1008801	-5.88	0.000	-.790831	-.3946519
Natl_WSU_EDF vs Natl_OTD_GTI	-.7789648	.108942	-7.15	0.000	-.9928849	-.5650447

## Linear and Logistic Regression of Concentration vs. Leak Flow Rates

### Linear Regressions

The linear and logistic regressions of maximum surface concentration vs. leak rate measurements were conducted to explore possible correlations, but were not used to inform the development of the Decision Tree process, including:

- The concentration thresholds,

- The average leak rate of the samples,
- The probability of the DT error type (i.e. the confusion/error matrix), and
- The calculation of the company-specific emission factors (EFs).

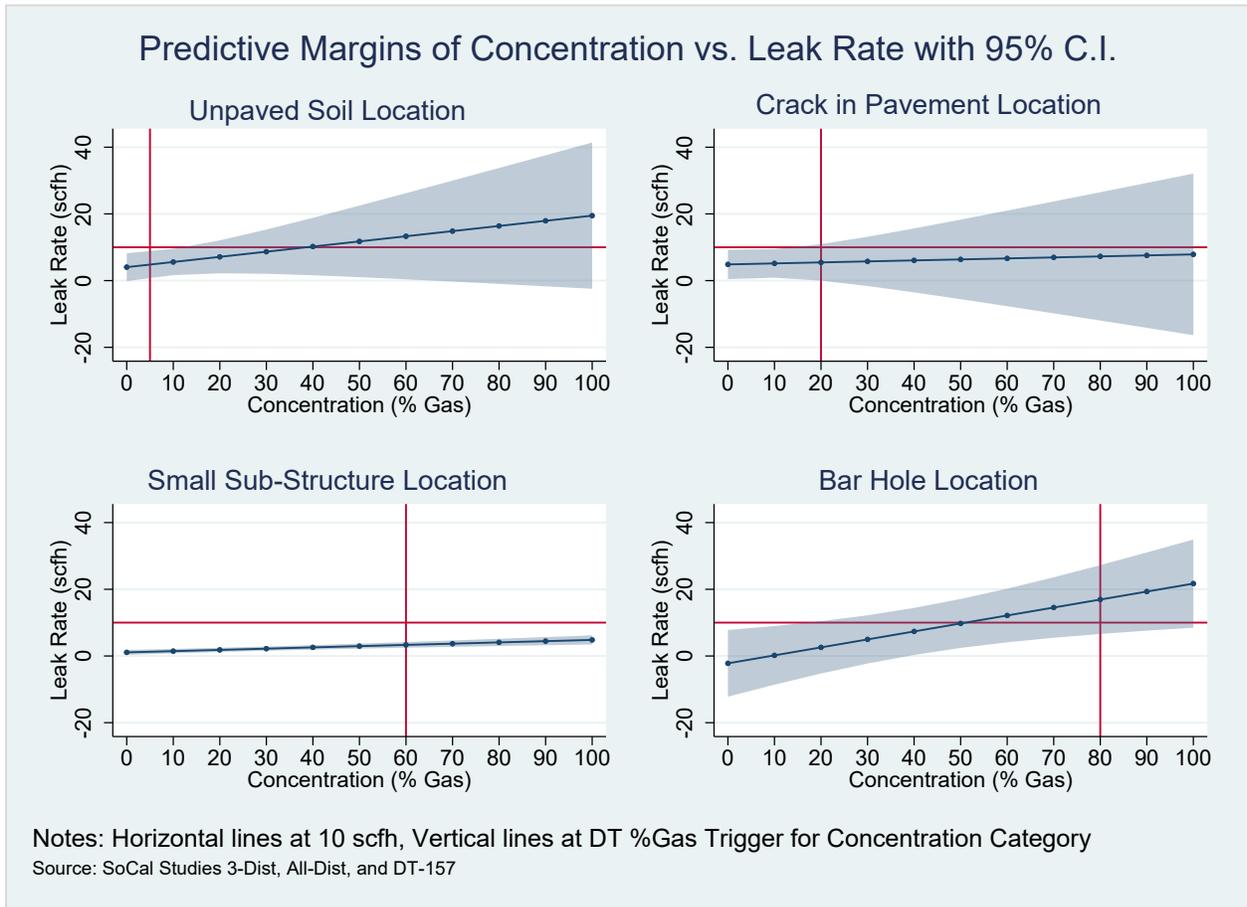
The two regression analyses of the relationship (or lack thereof) between the concentration and leak rate were in fact done post facto as a means of evaluating the empirical DT development. Since only the purely probabilistic Bayesian analysis was used to establish the important true and false negative and positive error values associated with the DT processes these regression sections are presented in the Appendix section of this report.

The regression margin plots in Figure 37 show average leak rate trend upward with increase in concentration, however, the 95% confidence intervals do cross zero in three of the cases and this would be expected from the regression results.

In the regression plot, the Decision Tree 10 scfh leak rate value that separates “Not Large” from “Large” non-hazardous leak levels is plotted as a *horizontal* line and each of the concentration levels that trigger a positive Decision Tree categorization are plotted as *vertical* lines at 80, 20, 60, and 5% gas respectively.

One can see that with the exception of the small sub-structure category, the intersection of the horizontal and vertical lines fall on or within the confidence intervals. The small sub-structure plot has tighter confidence intervals as well as a flatter regression slope for the average value. The trigger of 60% gas is conservatively set nonetheless, since the intersection of the lines is above the confidence interval bands. However, additional sampling and studies need to be conducted to increase the data set and improve the regression and associated confidence intervals.

**Figure 37: Predictive Margins for Leak Rate by Concentration from Linear Model.**



The linear regressions of concentration vs. leak rate are in Table 35 below. The regression was set up with the continuous (metric) independent variable as concentration and the dependent variable as leak flow rate. The regression was completed four times, once for each surface concentration category. The correlation is extremely poor, partly due to running the regression vs. the leak rate as opposed to the  $\log(10)$  of the leak rate. However, the analysis was done more for an illustrative purpose; to plot the result margins and show the expected value of the leak rate (not log if the same) vs. each type of surface concentration measure.

**Table 35: Linear Regression of DT Concentration and SoCalGas Study Leak Rates.**

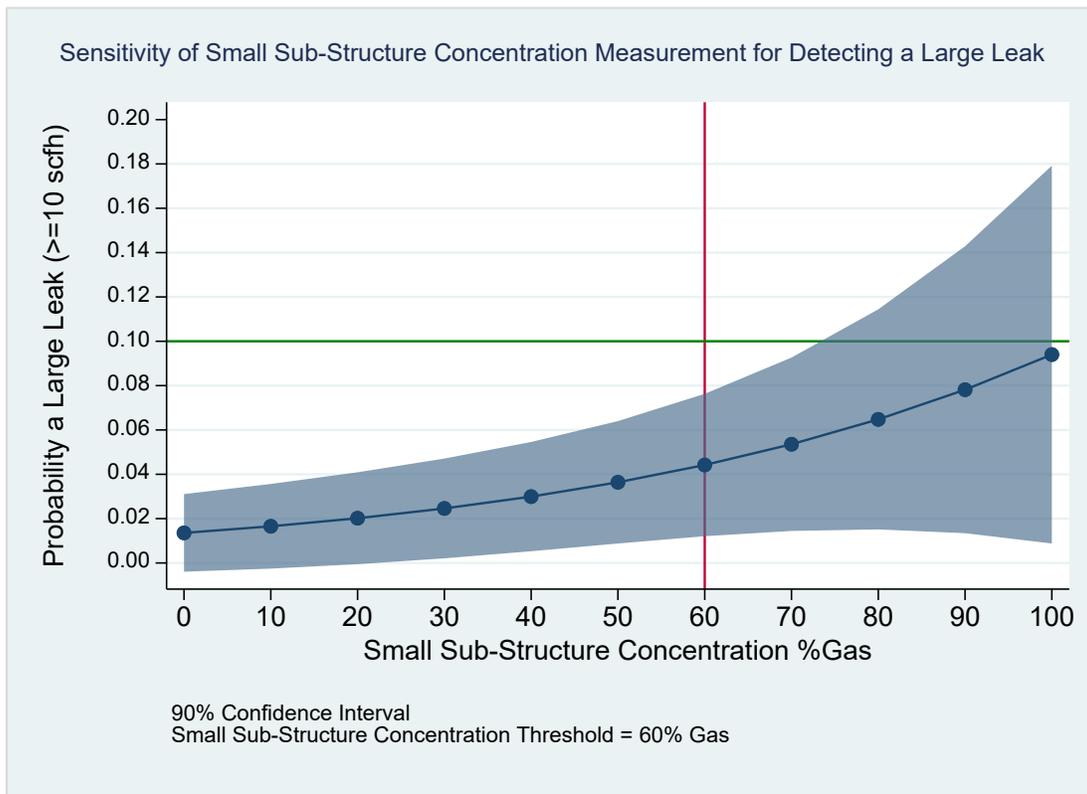
<b>Regress leakRate conc_bh_80</b>						
Source	SS	df	MS	Number of obs	=	
Model	9515.73041	1	9515.73041	F(1, 110)	=	6.72
Residual	155816.784	110	1416.51622	Prob > F	=	0.0108
Total	165332.514	111	1489.48211	R-squared	=	0.0576
				Adj R-squared	=	0.0490
				Root MSE	=	37.637
leakRate	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
conc_bh_80	.2391175	.0922574	2.59	0.011	.056285	.4219499
_cons	-2.20756	5.045601	-0.44	0.663	-12.20676	7.791637
<b>Regress leakRate conc_cip_20</b>						
Source	SS	df	MS	Number of obs	=	
Model	47.3922424	1	47.3922424	F(1, 194)	=	0.05
Residual	168596.222	194	869.052691	Prob > F	=	0.8156
Total	168643.614	195	864.839048	R-squared	=	0.0003
				Adj R-squared	=	-0.0049
				Root MSE	=	29.48
leakRate	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
conc_cip_20	.0300363	.1286223	0.23	0.816	-.2236413	.283714
_cons	4.863149	2.23957	2.17	0.031	.4461168	9.280181
<b>Regress leakRate conc_sss_60</b>						
Source	SS	df	MS	Number of obs	=	
Model	348.668852	1	348.668852	F(1, 164)	=	21.18
Residual	2700.05992	164	16.46378	Prob > F	=	0.0000
Total	3048.72878	165	18.4771441	R-squared	=	0.1144
				Adj R-squared	=	0.1090
				Root MSE	=	4.0576
leakRate	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
conc_sss_60	.0374123	.0081297	4.60	0.000	.02136	.0534647
_cons	1.077003	.3728651	2.89	0.004	.3407682	1.813238
<b>Regress leakRate conc_us_5</b>						
Source	SS	df	MS	Number of obs	=	
Model	1403.11046	1	1403.11046	F(1, 208)	=	1.73
Residual	168849.997	208	811.778832	Prob > F	=	0.1901
Total	170253.108	209	814.60817	R-squared	=	0.0082
				Adj R-squared	=	0.0035
				Root MSE	=	28.492
leakRate	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
conc_us_5	.154349	.1174023	1.31	0.190	-.077102	.3858001
_cons	4.044118	2.121006	1.91	0.058	-.1373072	8.225544

## Logistic Regressions

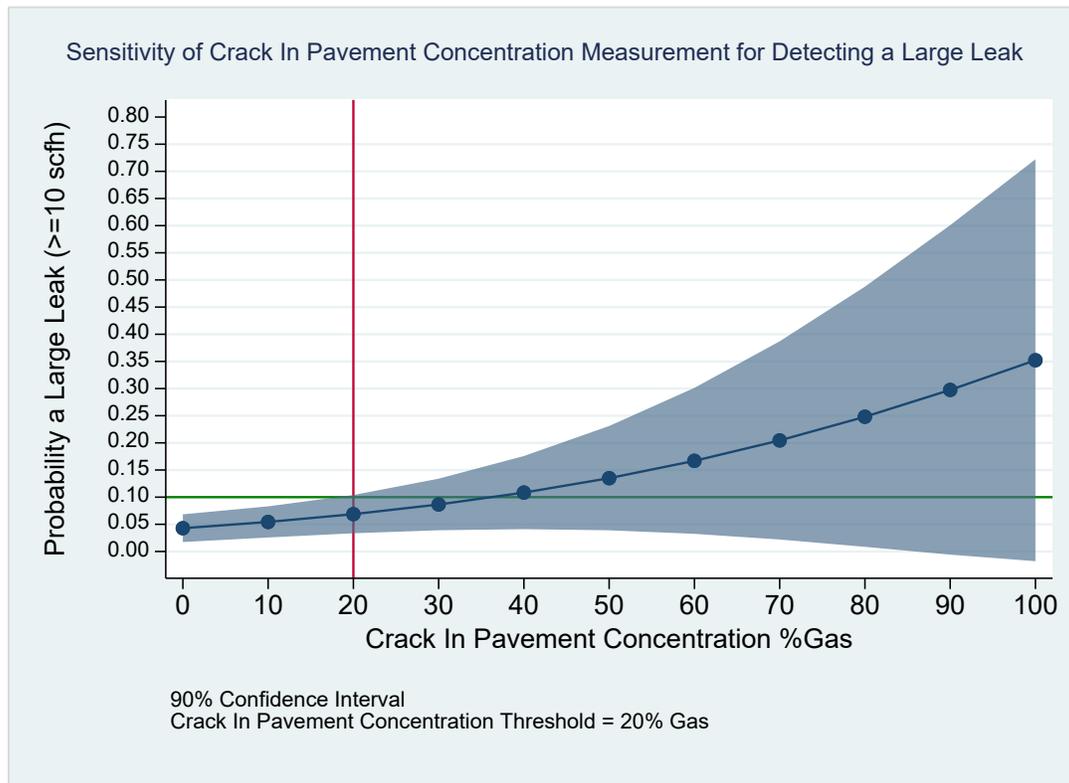
This section presents logistic regressions of the probability of a large leak rate as a function of concentration measurements. This analysis was done post facto and did not contribute to the DT thresholds for concentrations in the report. The basic steps for the regression include:

1. Establish the leak rate scfh large leak threshold. This decision can be based on several criteria, including the distribution of leak rate values from published industry studies and/or the company-specific leak rate distributions encountered in the field. For this study the value of 10 scfh was already established.
2. Collect field concentration and leak rate data as discussed in early sections of this report.
3. Run Logistic Regression with the continuous independent variable set to the leak concentration in %gas and the dependent categorical (binary) variable set to large (greater than or equal to 10 scfh) vs. not large (less than 10 scfh) leak rates. The logistic regression output margin plots for the SoCalGas DT concentration categories are presented in Figure 38 to Figure 40 below.

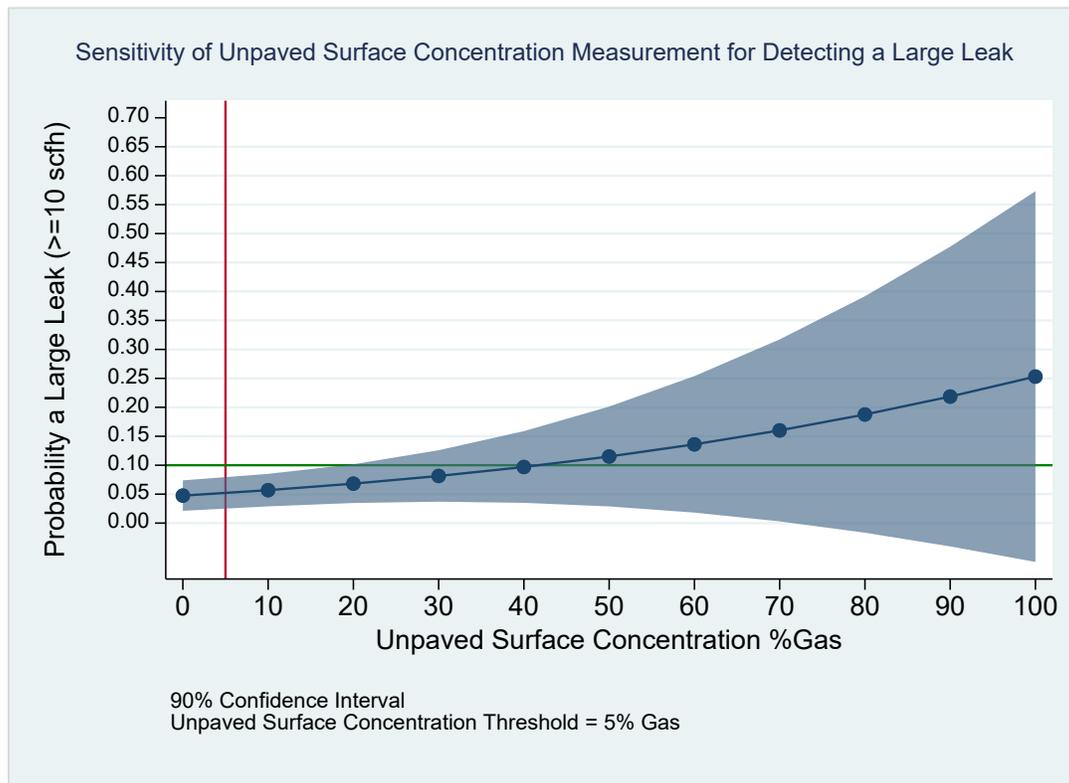
**Figure 38: Sensitivity of SSS Concentration to Large Leak Detection (Logistic Regression).**



**Figure 39: Sensitivity of CIP Concentration to Large Leak Detection (Logistic Regression).**

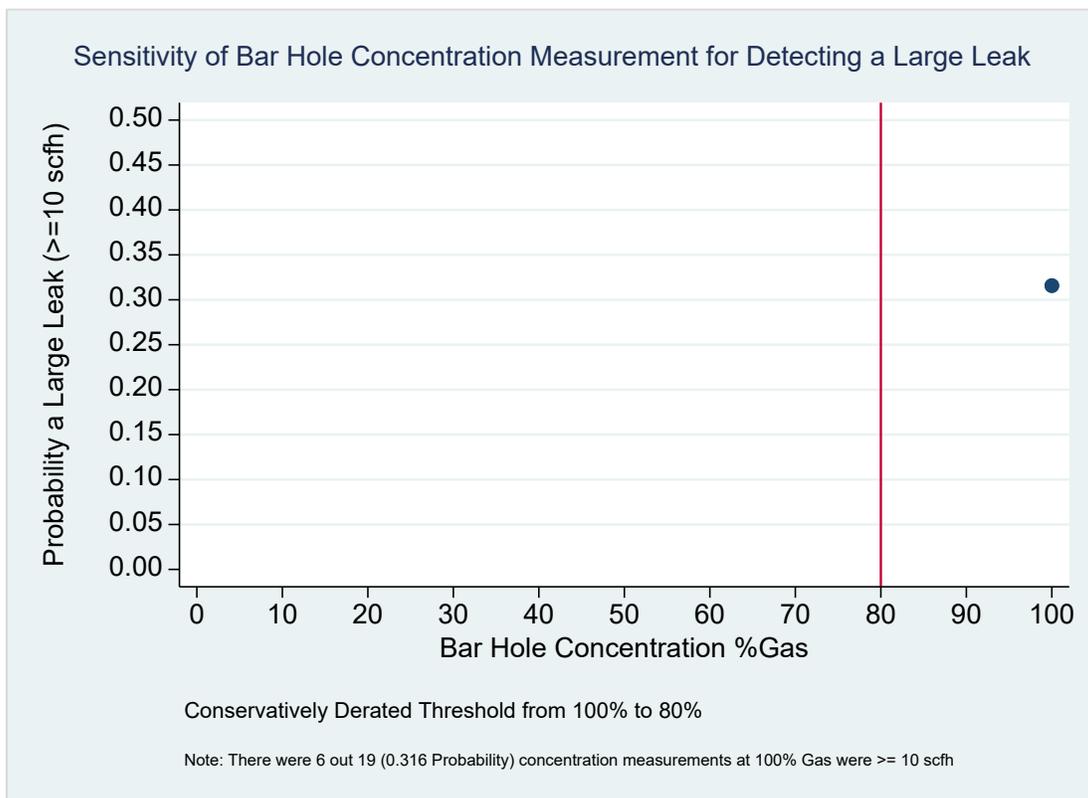


**Figure 40: Sensitivity of US Concentration to Large Leak Detection (Logistic Regression).**



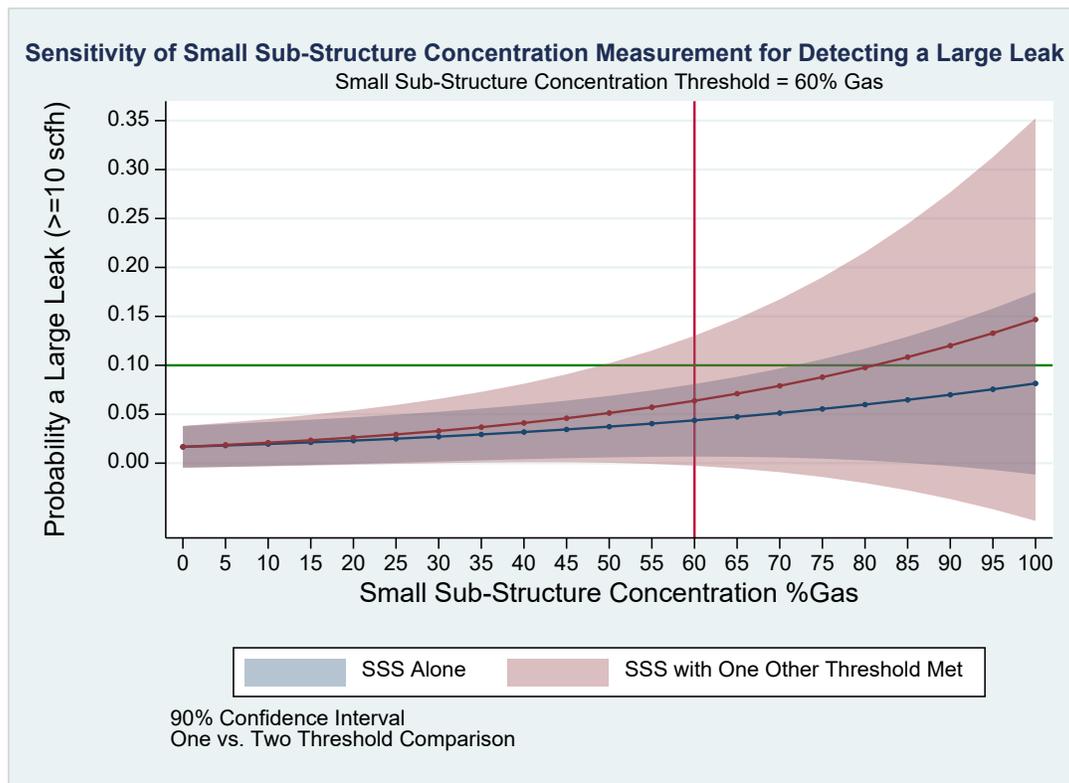
4. It is possible that the data will not support the *convergence* of a logistic (or other) regression. This was the case for the Bar Hole concentration vs. leak rate data as shown in Figure 41. To address these situations, and assign a concentration threshold, one should:
  - a. First, take a straight ratio of the relationship (as noted in the plot note). In this case, there was a 31.6% probability of a large leak when the Bar Hole concentration was 100% gas. This value was de-rated (lowered) by 20% from 100% to 80% for the threshold as a conservative measure.
  - b. Second, to establish the confidence in this measure, a Bayes analysis confirmed that if one has a gas concentration of less than or equal to 80% in a Bar Hole, then one should expect zero (0) probability of large leaks and will be 95% confident that the actual percent is no higher than 4.1%. This was based on having 71 data points (field samples) that had Bar Hole concentrations that were less than or equal to 80% gas and *all* 71 had less than or equal to 10 scfh measured leak rates.

**Figure 41: Sensitivity of BH Concentration to Large Leak Detection (Logistic Regression).**

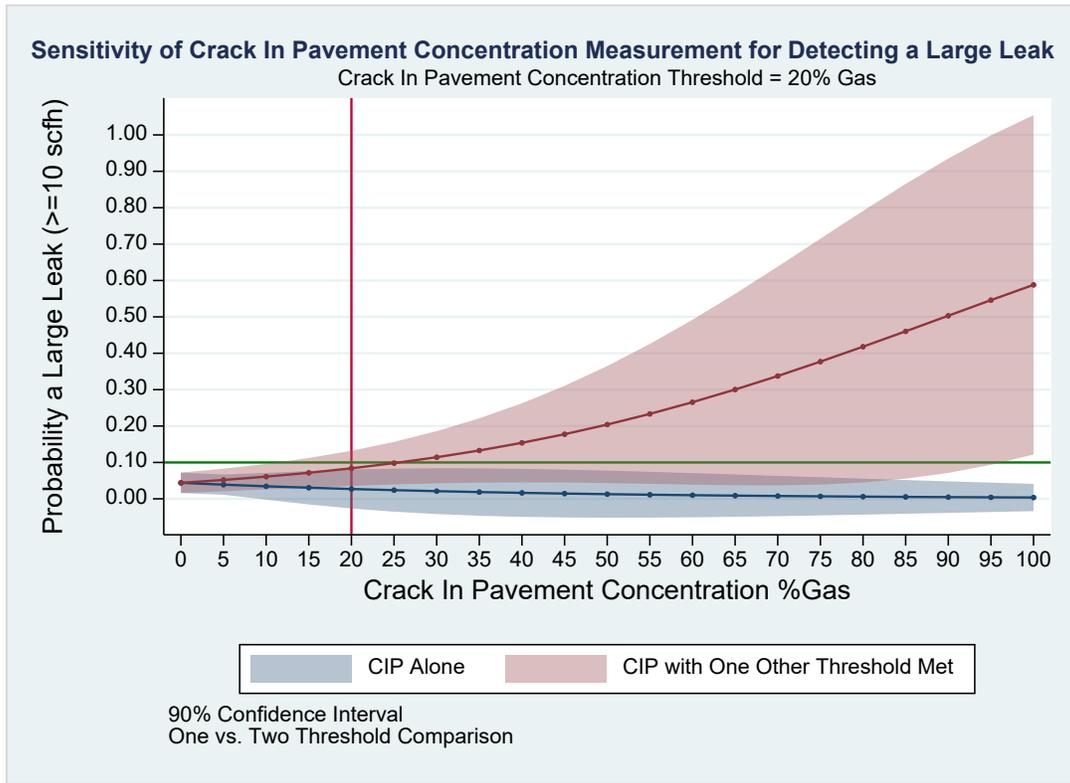


5. Establish the concentration thresholds for each category by setting the confidence interval to 90%, so you have a single sided upper limit of 95%. Then make an initial threshold assignment to establish a 95% confidence of a 10% probability (subjective decision on this number) or less of a large leak at or below the concentration threshold.
6. Then, with the new thresholds set, run a second Logistic Regression with *interactions* this time to take credit for more than one threshold when it is triggered. See Figure 42 to Figure 44 below.

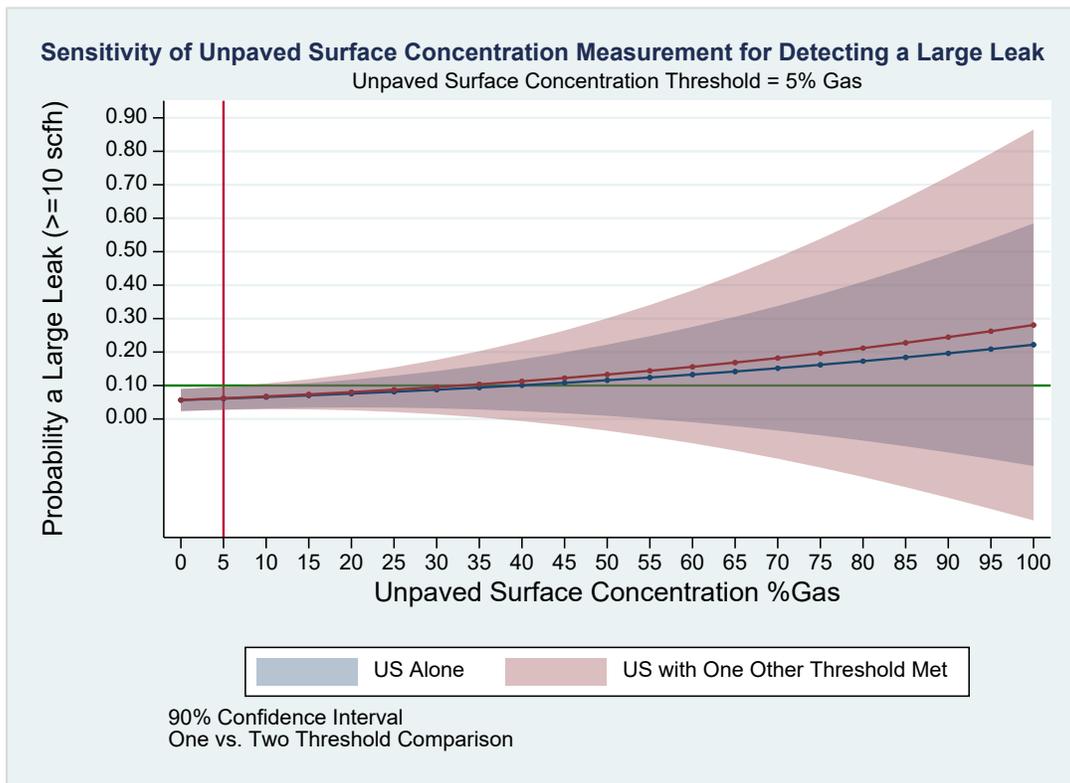
**Figure 42: Sensitivity of SSS Concentration to Large Leak Detection - Multi-Thresholds.**



**Figure 43: Sensitivity of CIP Concentration to Large Leak Detection - Multi-Thresholds.**



**Figure 44: Sensitivity of US Concentration to Large Leak Detection - Multi-Thresholds.**



7. Check that (a) the single thresholds *still* have the 95% confidence of no more than a 10% chance of large leak (which should be the case); but now you can (b) check that when you have two thresholds met, that the expected (mean) value for that is also below the 10% probability. The concentration threshold can be adjusted to meet both these objectives.
8. Based on this analysis, for the SoCalGas study, one could change the SSS threshold from 60% gas to 70% gas, the other three are set as noted in (7) above. The 60% gas SSS concentration threshold is more conservative than 70% gas and reflects that a 60% gas concentration is considered "hazardous" for this category in California.
9. If the SSS threshold was increased from 60% gas to 70% gas, it would have improved the DT results. There would be four (4) less false positives and no other changes, i.e. no increases in false negatives across the 291 sample data set.

# Appendix C: Log-normal Distribution Facts

## Log-normal Distribution Equations

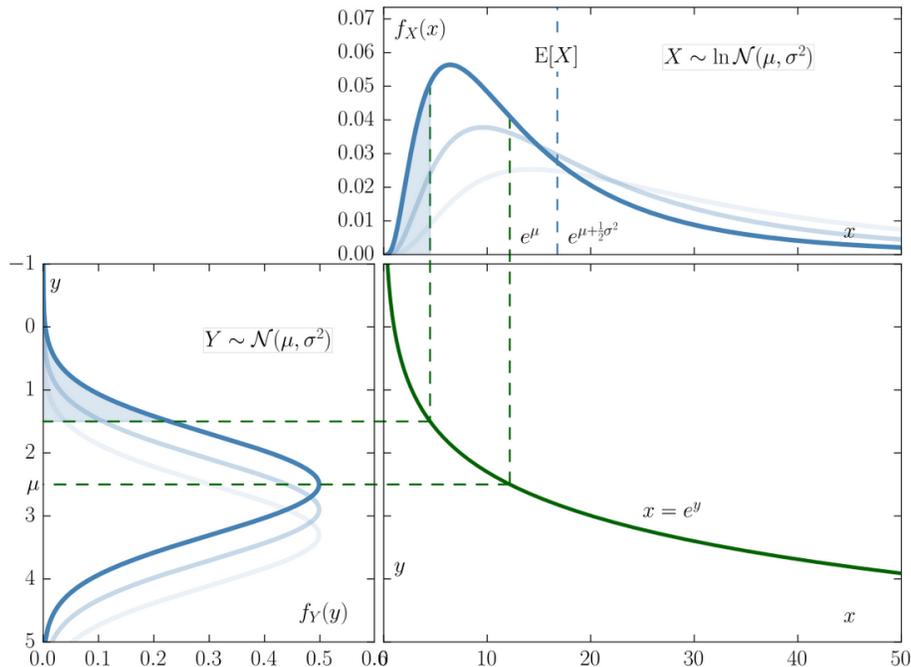
Table 36: Log-Normal Distribution Equations.

Probability density function :	$f(x) = \frac{1}{x\sqrt{2\pi\sigma_1^2}} \exp\left[-\frac{(\ln[x] - \mu_1^2)}{2\sigma_1^2}\right]$
	<p>where <math>\mu_1 = \ln\left[\frac{\mu^2}{\sqrt{\sigma^2 + \mu^2}}\right]</math> and <math>\sigma_1 = \sqrt{\ln\left[\frac{\sigma^2 + \mu^2}{\mu^2}\right]}</math></p>
Cumulative distribution function :	No closed form
Parameter restriction :	$\sigma > 0, \mu > 0$
Domain :	$x \geq 0$
Mean :	$\mu$
Mode :	$\exp(\mu_1 - \sigma_1^2)$
Variance :	$\sigma^2$
Skewness :	$\left(\frac{\sigma}{\mu}\right)^3 + 3\left(\frac{\sigma}{\mu}\right)$
Kurtosis :	$z^4 + 2z^3 + 3z^2 - 3$ where $z = 1 + \frac{\sigma}{\mu}$

Source: [34]

## Log-normal Plot and Comparison to Normal Distribution

Figure 45: Log-normal Distribution Theoretical Plot.



Source: [35]

## Distribution Fitting Goodness of Fit

Although still popular today, the Chi-Squared, Kolmogorov-Smirnoff and Anderson-Darling goodness of fit statistics are technically all inappropriate as a method of comparing fits of distributions to data[36].

They are also limited to having precise observations and cannot incorporate censored, truncated or binned data. Realistically, most of the time we are fitting a continuous distribution to a set of precise observations and then the Anderson-Darling does a reasonable job.

For important work you should instead consider using statistical measures of fit called information criteria.

- **SIC** (Schwarz information criterion, aka Bayesian information criterion BIC)[37]
- **AIC** (Akaike information criterion)[38]
- **HQIC** (Hannan-Quinn information criterion)[39]+

The aim is to find the model with the lowest value of the selected information criterion.

We decided to illustrate the log-normal distribution, but as can be seen in Table 37 the log-gamma has an excellent fit. The log-normal was used for illustration due to its common application in these types of log distributions.

**Table 37: Distribution Goodness of Fit to SoCalGas data set (291 samples).**

Distributions fitted	Data	LogGamma	Dagum	Log-normal	LogLogistic	LogLaplace	Weibull
<i>Goodness of fit</i>							
AIC		966.8	983.9	990.1	991.2	994.2	1014.7
AIC ranking		1	2	3	4	5	6
SIC		977.7	994.9	997.4	998.5	1005.1	1022.0
SIC ranking		1	2	3	4	5	6
HQIC		971.1	988.3	993.0	994.1	998.5	1017.6
HQIC ranking		1	2	3	4	5	6
<i>Comparison of data and fitted distribution statistics</i>							
Minimum	0.00272	0.00272	0	0	0	0	0
Maximum	373	+Infinity	+Infinity	+Infinity	+Infinity	+Infinity	+Infinity
Mean	4.31	Undefined	4.99	4.67	Undefined	Undefined	3.14
St. Dev	24.37	Undefined	Undefined	32.30	Undefined	Undefined	6.50

# Appendix D: Leak Spread Comparison to Leak Rate

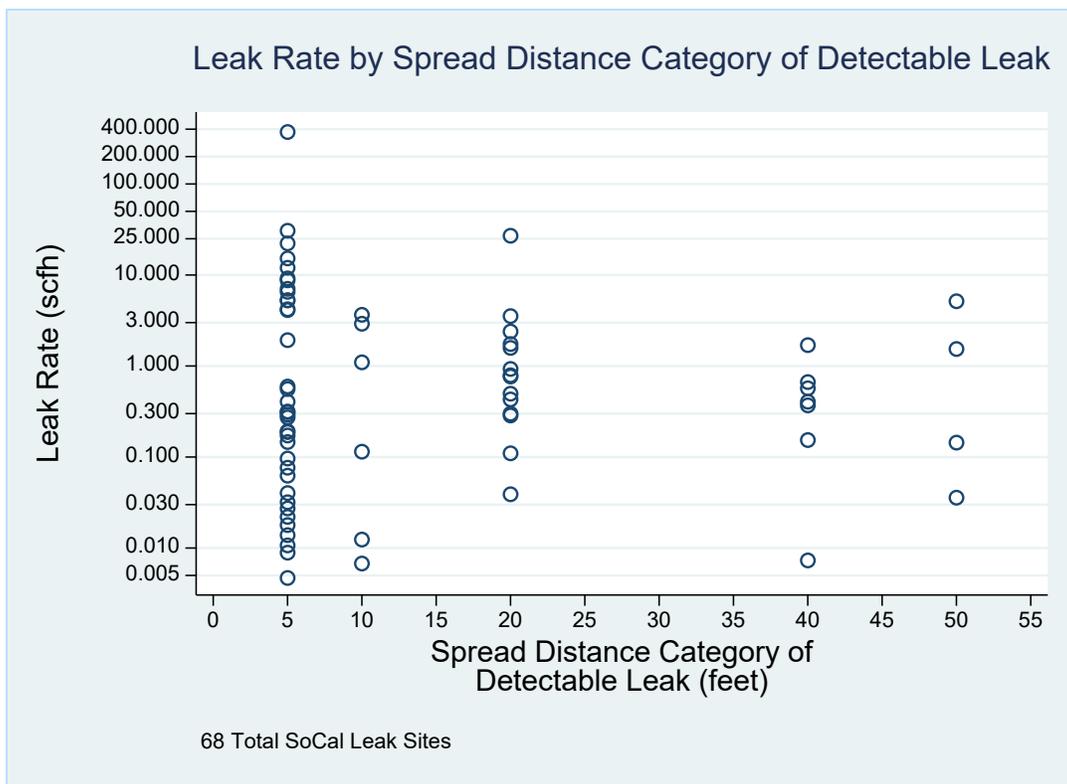
As noted in the Background section of this report, early in the development of the program, system data were mined, and system leak flow rates were measured to evaluate the relationships between measured methane concentration and “leak spread” data against the flow rate of the leak. However, no correlation was found in this data, see Figure 46 to Figure 49.

A total of 68 leak sites at SoCalGas had measured methane concentration levels and the associated leak spread. The largest spread (distance in feet) across the leak site was recorded and placed into one of four categories by length range, see Table 38. These four categories were assigned a Numeric Code as shown in the table. The numeric code is used in the figures as well.

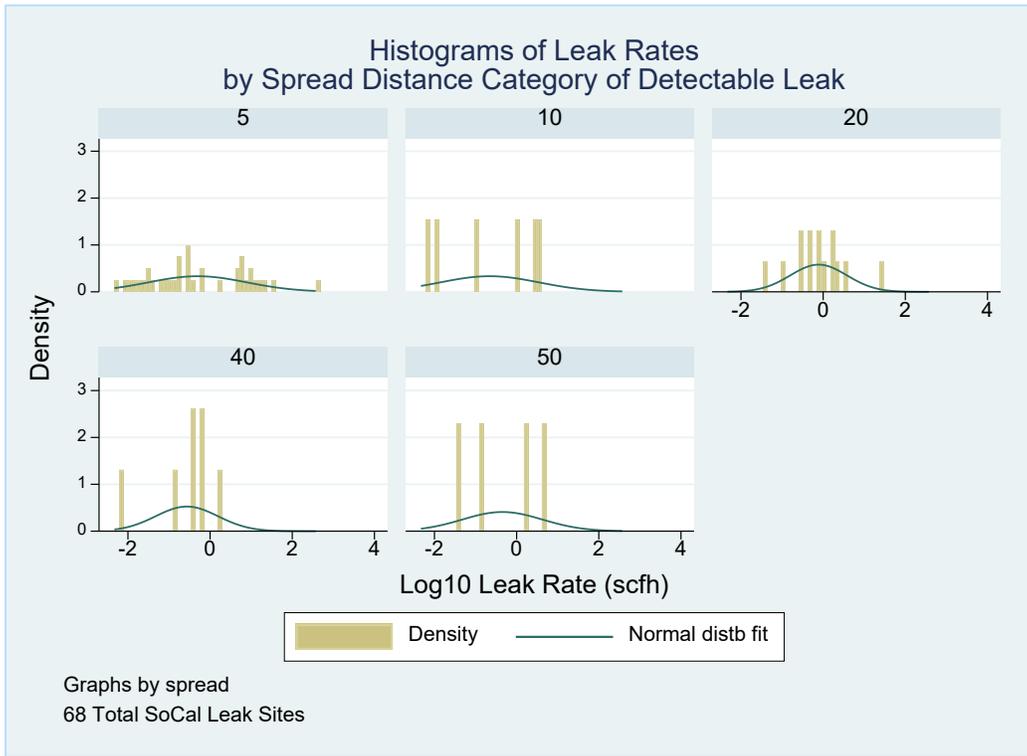
**Table 38: Leak Spread Categories.**

Category Length Range (feet)	Category Numeric Code
0 to 5	5
6 to 10	10
11 to 20	20
21 to 40	40
> 40 feet	50

**Figure 46: Scatter Plot of Leak Concentration Spread vs. Leak Rate.**



**Figure 47: Histogram of Leak Rates by Leak Concentration Spread.**



**Figure 48: Box Plot of Leak Concentration Spread vs. Leak Rate.**

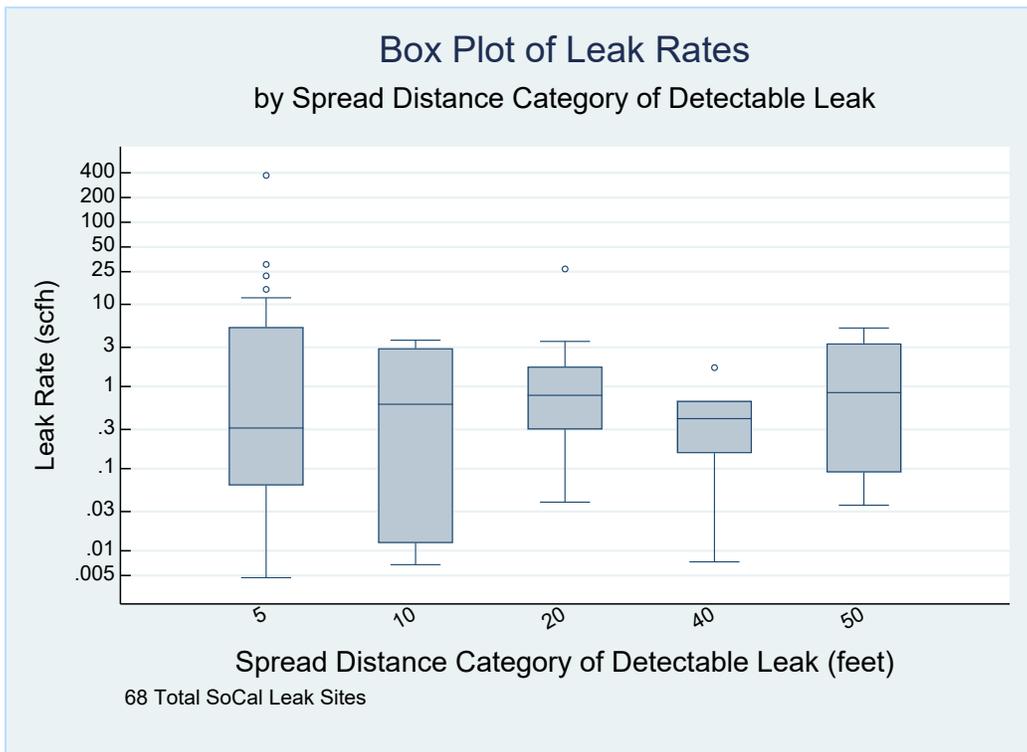
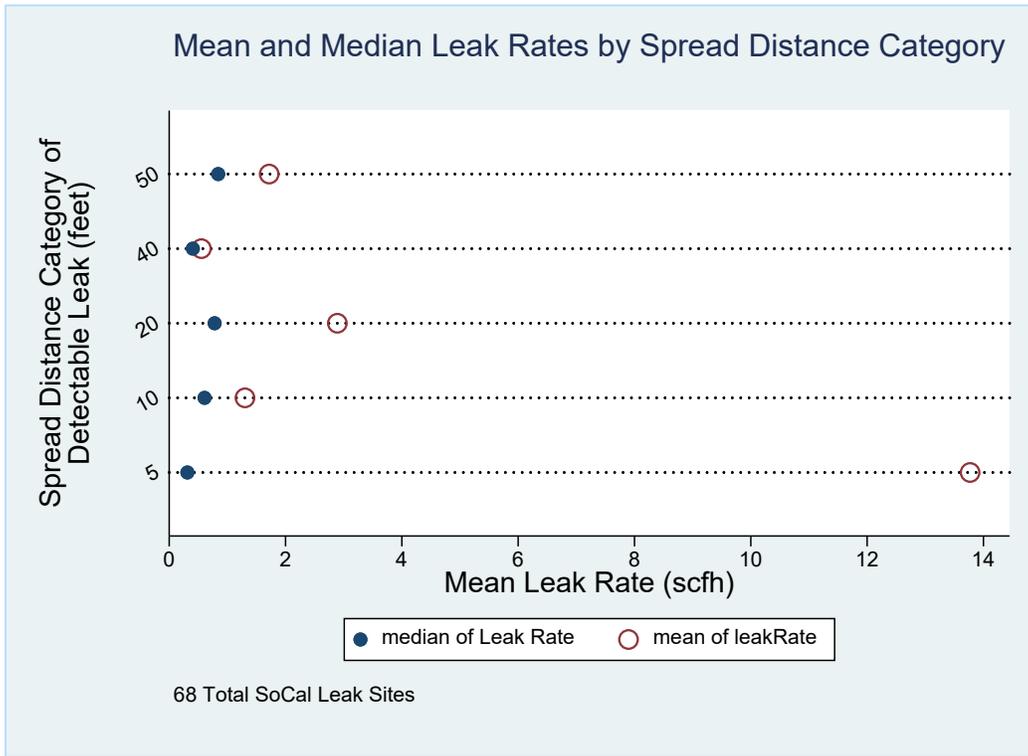


Figure 49: Mean and Median of Leak Rate by Leak Concentration Spread.



# Appendix E: Study Leak Rate and Concentration Data

For the SoCalGas studies, the DT approach collected methane concentration measurements at defined types of surface condition locations. The four defined types of surface condition locations are listed in the Table 39 below and use the following abbreviations:

- Bar Hole (leak survey type) - BH
- Crack (or seam) In Pavement - CIP
- Small Sub-Structure (not gas system related) - SSS
- Unpaved Surface - US

**Table 39: Leak Rate scfh Methane (CH<sub>4</sub>) and Concentration (% gas) by Study.**

	Study	Leak Rate (scfh)	Conc. BH (% gas)	Conc. CIP (% gas)	Conc. SSS (% gas)	Conc. US (% gas)
1.	AllDisPilot	0.0027	.	0.00	0.10	0.00
2.	AllDisPilot	0.0047	.	0.00	.	0.00
3.	AllDisPilot	0.0055	.	0.00	0.29	0.00
4.	AllDisPilot	0.0068	.	0.17	.	0.00
5.	AllDisPilot	0.0073	0.01	0.00	0.40	0.00
6.	AllDisPilot	0.0089	.	0.07	0.00	0.05
7.	AllDisPilot	0.0107	.	0.00	.	0.01
8.	AllDisPilot	0.0124	.	0.01	0.00	0.00
9.	AllDisPilot	0.0136	.	0.23	0.00	0.00
10.	AllDisPilot	0.0139	0.00	0.00	0.01	0.00
11.	AllDisPilot	0.0179	0.50	0.00	0.00	0.00
12.	AllDisPilot	0.0199	.	0.00	0.07	0.00
13.	AllDisPilot	0.0273	.	0.02	0.00	0.01
14.	AllDisPilot	0.0300	0.00	0.03	0.00	0.00
15.	AllDisPilot	0.0319	.	1.60	0.04	0.00
16.	AllDisPilot	0.0358	0.00	0.00	.	0.30
17.	AllDisPilot	0.0390	0.00	0.00	0.00	0.04
18.	AllDisPilot	0.0403	0.15	0.00	0.00	0.00
19.	AllDisPilot	0.0409	0.02	0.00	0.00	0.00
20.	AllDisPilot	0.0575	.	0.05	.	0.50
21.	AllDisPilot	0.0614	.	0.00	0.50	0.00
22.	AllDisPilot	0.0625	0.06	0.00	0.00	0.00
23.	AllDisPilot	0.0765	.	0.00	0.00	0.20
24.	AllDisPilot	0.0769	0.00	0.00	0.70	0.00
25.	AllDisPilot	0.0798	0.00	0.00	0.00	0.05
26.	AllDisPilot	0.0964	.	0.20	0.00	.
27.	AllDisPilot	0.1101	.	0.00	0.02	0.00
28.	AllDisPilot	0.1147	1.20	0.00	0.01	0.18
29.	AllDisPilot	0.1329	.	0.00	0.00	0.04
30.	AllDisPilot	0.1442	.	0.00	0.06	0.00
31.	AllDisPilot	0.1465	.	0.28	0.00	0.00
32.	AllDisPilot	0.1540	.	0.00	0.23	0.00
33.	AllDisPilot	0.1726	.	0.00	0.00	0.13
34.	AllDisPilot	0.1884	.	0.00	0.00	0.00
35.	AllDisPilot	0.1894	.	0.06	0.20	0.11
36.	AllDisPilot	0.1991	25.00	0.01	.	0.00
37.	AllDisPilot	0.2168	.	0.00	0.00	0.26
38.	AllDisPilot	0.2263	.	0.00	0.00	0.04
39.	AllDisPilot	0.2367	.	0.00	0.00	0.06

	Study	Leak Rate (scfh)	Conc. BH (% gas)	Conc. CIP (% gas)	Conc. SSS (% gas)	Conc. US (% gas)
40.	AllDisPilot	0.2710	.	0.20	0.00	0.20
41.	AllDisPilot	0.2730	.	0.00	.	0.03
42.	AllDisPilot	0.2749	0.00	0.00	0.00	1.30
43.	AllDisPilot	0.2823	.	0.15	1.60	0.01
44.	AllDisPilot	0.2867	42.00	0.02	0.00	0.00
45.	AllDisPilot	0.2920	.	0.30	0.00	0.01
46.	AllDisPilot	0.3120	2.00	0.00	0.00	0.01
47.	AllDisPilot	0.3168	.	0.30	0.20	0.72
48.	AllDisPilot	0.3319	.	1.00	0.00	0.02
49.	AllDisPilot	0.3879	.	0.20	0.00	0.00
50.	AllDisPilot	0.4060	.	0.01	0.00	3.50
51.	AllDisPilot	0.4060	0.85	0.06	0.00	0.06
52.	AllDisPilot	0.4975	.	0.00	0.00	0.30
53.	AllDisPilot	0.5600	.	0.00	0.00	0.34
54.	AllDisPilot	0.5698	.	0.01	0.06	0.00
55.	AllDisPilot	0.6676	.	3.00	0.00	.
56.	AllDisPilot	0.7387	.	0.03	0.01	0.20
57.	AllDisPilot	0.7688	.	7.50	0.00	0.00
58.	AllDisPilot	0.7927	.	0.10	.	0.00
59.	AllDisPilot	1.3139	.	0.00	0.00	1.00
60.	AllDisPilot	1.5844	.	0.02	0.00	0.97
61.	AllDisPilot	1.6990	.	0.00	10.00	0.00
62.	AllDisPilot	1.7470	.	0.00	0.20	0.13
63.	AllDisPilot	1.9275	.	1.60	0.00	0.17
64.	AllDisPilot	2.2184	0.00	0.00	0.00	0.45
65.	AllDisPilot	2.3995	.	2.95	0.00	0.01
66.	AllDisPilot	2.9205	0.00	0.00	0.18	0.35
67.	AllDisPilot	4.1216	0.01	0.00	.	0.05
68.	AllDisPilot	4.2127	.	1.00	4.00	0.10
69.	AllDisPilot	5.1560	.	1.00	.	4.00
70.	AllDisPilot	5.2590	.	0.00	1.20	0.00
71.	AllDisPilot	5.2750	.	0.00	0.00	2.00
72.	AllDisPilot	7.3507	23.00	0.00	.	0.70
73.	AllDisPilot	9.1915	.	0.10	0.10	1.45
74.	AllDisPilot	12.0244	.	8.00	.	1.00
75.	AllDisPilot	15.2524	.	0.00	0.80	0.01
76.	AllDisPilot	27.0447	.	40.00	.	20.00
77.	AllDisPilot	0.0220	.	0.09	0.00	0.00
78.	AllDisPilot	0.5958	.	0.00	40.00	0.00
79.	AllDisLIRP	0.4301	0.00	0.00	5.00	0.00
80.	AllDisLIRP	0.3704	.	0.00	.	0.04
81.	AllDisLIRP	30.7021	.	0.00	.	6.00
82.	AllDisLIRP	0.1922	38.00	0.00	.	0.20
83.	AllDisLIRP	3.6600	.	4.60	.	0.33
84.	AllDisLIRP	3.5400	100.00	1.60	.	0.90
85.	AllDisLIRP	1.5400	19.00	0.50	.	52.00
86.	AllDisLIRP	1.1000	95.00	0.01	0.01	0.13
87.	AllDisLIRP	0.9300	75.00	1.80	0.13	0.06
88.	AllDisLIRP	0.3000	0.39	0.00	0.00	0.12
89.	3DisPilotLowSpec	0.0800	.	3.00	.	3.00
90.	3DisPilotLowSpec	0.1100	70.00	1.00	.	1.00
91.	3DisPilotLowSpec	1.6400	20.00	.	5.00	3.00
92.	3DisPilotLowSpec	0.3900	1.00	0.00	0.20	1.00
93.	3DisPilotLowSpec	0.0600	30.00	.	0.30	1.00
94.	3DisPilotLowSpec	0.0600	4.00	0.05	0.00	4.00
95.	3DisPilotLowSpec	0.9000	40.00	0.00	0.00	2.00
96.	3DisPilotLowSpec	0.3400	40.00	2.00	0.00	2.00
97.	3DisPilot	373.0000	100.00	0.08	.	20.00
98.	3DisPilot	0.1200	70.00	50.00	7.00	25.00
99.	3DisPilot	0.3700	20.00	15.00	0.03	5.00

	Study	Leak Rate (scfh)	Conc. BH (% gas)	Conc. CIP (% gas)	Conc. SSS (% gas)	Conc. US (% gas)
100.	3DisPilot	0.2700	.	8.00	.	31.00
101.	3DisPilot	10.0000	100.00	1.00	.	2.00
102.	3DisPilot	22.2900	100.00	0.70	.	5.00
103.	3DisPilot	4.0600	100.00	4.00	0.00	0.00
104.	3DisPilot	1.2400	100.00	.	.	15.00
105.	3DisPilot	0.4800	35.00	35.00	0.00	.
106.	3DisPilot	0.9600	.	.	64.00	.
107.	3DisPilot	0.4400	92.00	1.00	0.00	1.00
108.	3DisPilot	0.1100	0.00	0.00	0.00	70.00
109.	3DisPilot	0.2100	85.00	0.15	0.33	.
110.	3DisPilot	0.3300	70.00	0.00	0.00	10.00
111.	3DisPilot	0.0830	52.00	0.00	0.00	26.00
112.	3DisPilot	0.9100	85.00	.	.	.
113.	3DisPilot	7.0600	85.00	0.02	0.00	0.44
114.	3DisPilot	6.5500	10.00	0.04	.	10.00
115.	3DisPilot	5.2900	5.00	0.03	.	5.00
116.	3DisPilot	1.1100	67.00	1.50	.	9.50
117.	3DisPilot	8.7500	100.00	0.02	0.00	0.08
118.	3DisPilot	0.1115	89.00	0.01	.	1.00
119.	3DisPilot	2.2960	27.00	0.10	.	23.00
120.	3DisPilot	0.1100	80.00	0.00	0.00	80.00
121.	3DisPilot	172.4400	100.00	7.00	.	3.00
122.	3DisPilot	0.0500	10.00	1.50	0.00	0.00
123.	3DisPilot	3.2000	7.00	0.40	.	0.02
124.	3DisPilot	0.6100	25.00	0.10	.	0.10
125.	3DisPilot	0.2400	22.00	0.01	.	0.60
126.	3DisPilot	0.0300	19.00	0.02	0.00	0.38
127.	3DisPilot	0.0300	2.00	0.05	0.61	0.15
128.	3DisPilot	0.9200	48.00	0.01	0.00	0.02
129.	3DisPilot	3.1200	62.00	0.30	.	0.10
130.	3DisPilot	0.0400	0.24	0.05	.	0.02
131.	3DisPilot	0.6100	0.00	0.00	10.00	0.00
132.	3DisPilot	3.8200	0.20	.	.	0.05
133.	3DisPilot	1.8400	25.00	0.09	.	0.35
134.	3DisPilot	5.5400	0.00	0.00	0.00	0.01
135.	3DisPilot	0.5700	30.00	0.01	0.00	0.11
136.	3DisPilot	0.0200	4.00	0.00	0.00	0.14
137.	3DisPilot	0.3300	.	0.10	0.00	0.00
138.	3DisPilot	0.0600	0.09	0.03	0.00	0.02
139.	3DisPilot	2.1800	100.00	1.50	.	2.50
140.	3DisPilot	2.2300	100.00	6.00	.	0.00
141.	3DisPilot	3.6700	100.00	5.00	1.00	2.00
142.	3DisPilot	6.9300	42.00	0.80	0.00	9.00
143.	3DisPilot	1.1000	25.00	0.02	0.40	12.00
144.	3DisPilot	0.0200	.	0.00	0.03	.
145.	3DisPilot	0.6900	0.00	0.00	100.00	0.00
146.	3DisPilot	0.1400	80.00	0.00	0.00	0.00
147.	3DisPilot	0.9800	25.00	0.03	3.00	0.03
148.	3DisPilot	4.5300	30.00	0.00	0.00	30.00
149.	3DisPilot	4.4400	12.00	0.01	.	9.00
150.	3DisPilot	1.0200	17.00	0.00	0.17	12.00
151.	3DisPilot	1.2400	20.00	1.00	0.00	12.00
152.	3DisPilot	0.5300	0.00	0.00	7.50	0.00
153.	DT157Pilot	0.0030	.	0.03	.	.
154.	DT157Pilot	0.0337	.	.	.	.
155.	DT157Pilot	1.1364	.	.	.	0.20
156.	DT157Pilot	0.6868	.	.	.	.
157.	DT157Pilot	0.3201	.	.	.	.
158.	DT157Pilot	0.0456	.	.	0.40	.
159.	DT157Pilot	0.0857	0.50	.	.	.

	Study	Leak Rate (scfh)	Conc. BH (% gas)	Conc. CIP (% gas)	Conc. SSS (% gas)	Conc. US (% gas)
160.	DT157Pilot	0.3121	.	0.61	.	.
161.	DT157Pilot	0.2598	0.90	.	.	.
162.	DT157Pilot	0.0471	.	1.08	.	0.08
163.	DT157Pilot	2.6745	.	.	.	1.10
164.	DT157Pilot	0.3009	.	.	1.10	.
165.	DT157Pilot	0.0420	.	1.12	.	.
166.	DT157Pilot	1.3500	.	.	.	1.14
167.	DT157Pilot	0.3674	.	.	.	.
168.	DT157Pilot	0.8960	.	.	52.00	.
169.	DT157Pilot	2.1040	.	.	.	1.50
170.	DT157Pilot	0.4980	.	.	.	1.50
171.	DT157Pilot	0.4934	.	.	.	.
172.	DT157Pilot	0.1844	24.00	.	.	.
173.	DT157Pilot	1.7206	.	2.50	.	.
174.	DT157Pilot	1.8764	.	.	3.50	.
175.	DT157Pilot	0.9169	.	3.70	.	.
176.	DT157Pilot	0.3302	.	6.00	.	.
177.	DT157Pilot	0.7801	.	.	7.00	.
178.	DT157Pilot	1.5748	.	.	.	.
179.	DT157Pilot	2.0911	.	.	8.00	.
180.	DT157Pilot	0.1558	.	.	32.00	.
181.	DT157Pilot	6.1970	.	.	.	8.50
182.	DT157Pilot	1.2619	.	.	.	.
183.	DT157Pilot	6.6519	.	3.00	.	12.00
184.	DT157Pilot	4.7489	.	13.00	.	.
185.	DT157Pilot	3.2550	.	5.00	16.00	1.50
186.	DT157Pilot	0.2394	.	2.00	.	16.00
187.	DT157Pilot	0.4809	18.00	.	.	.
188.	DT157Pilot	1.0130	20.00	.	.	.
189.	DT157Pilot	3.7304	20.00	.	.	.
190.	DT157Pilot	1.6414	23.00	.	.	.
191.	DT157Pilot	0.3124	25.00	.	.	.
192.	DT157Pilot	0.0444	.	.	32.00	.
193.	DT157Pilot	0.0456	.	.	13.00	.
194.	DT157Pilot	0.6210	.	.	.	.
195.	DT157Pilot	2.5850	.	.	50.00	.
196.	DT157Pilot	1.7481	52.00	.	22.00	0.50
197.	DT157Pilot	1.4472	.	.	59.00	.
198.	DT157Pilot	2.1470	.	.	100.00	.
199.	DT157Pilot	5.9664	100.00	0.50	.	.
200.	DT157Pilot	1.1704	.	.	100.00	.
201.	DT157Pilot	12.2475	.	2.00	100.00	.
202.	DT157Pilot	0.9859	.	.	100.00	.
203.	DT157Pilot	1.8368	.	.	100.00	.
204.	DT157Pilot	3.6167	100.00	.	100.00	.
205.	DT157Pilot	7.5300	.	23.00	100.00	.
206.	DT157Pilot	18.2144	.	40.00	100.00	2.00
207.	DT157Pilot	4.0392	.	100.00	100.00	.
208.	DT157Pilot	4.3662	.	.	100.00	10.00
209.	DT157Pilot	0.0329	1.00	.	.	.
210.	DT157Pilot	0.4608	.	.	.	1.00
211.	DT157Pilot	0.5173	.	.	.	1.70
212.	DT157Pilot	0.8978	.	.	.	9.00
213.	DT157Pilot	0.1640	.	1.60	.	.
214.	DT157Pilot	0.0158	.	0.15	.	.
215.	DT157Pilot	0.0173	.	.	.	0.01
216.	DT157Pilot	0.3845	.	3.00	.	.
217.	DT157Pilot	1.4855	.	10.00	.	.
218.	DT157Pilot	3.5326	.	.	.	5.60
219.	DT157Pilot	0.1560	.	.	.	0.12

	Study	Leak Rate (scfh)	Conc. BH (% gas)	Conc. CIP (% gas)	Conc. SSS (% gas)	Conc. US (% gas)
220.	DT157Pilot	0.2646	.	.	.	1.00
221.	DT157Pilot	0.0342	.	.	100.00	.
222.	DT157Pilot	20.0870	100.00	100.00	.	.
223.	DT157Pilot	1.4640	.	.	.	2.50
224.	DT157Pilot	2.9820	.	.	100.00	.
225.	DT157Pilot	5.1899	.	.	.	9.00
226.	DT157Pilot	0.4406	.	40.00	.	0.14
227.	DT157Pilot	1.2760	.	75.00	.	.
228.	DT157Pilot	2.5837	.	1.30	.	.
229.	DT157Pilot	1.1464	.	1.45	.	1.23
230.	DT157Pilot	1.0416	70.00	.	.	.
231.	DT157Pilot	3.9882	.	80.00	.	100.00
232.	DT157Pilot	0.7686	.	.	100.00	.
233.	DT157Pilot	1.5258	.	.	.	1.70
234.	DT157Pilot	2.0598	.	20.00	100.00	.
235.	DT157Pilot	3.0195	.	8.00	100.00	.
236.	DT157Pilot	1.9665	.	100.00	.	.
237.	DT157Pilot	6.2496	.	3.00	.	8.00
238.	DT157Pilot	0.7680	.	.	100.00	.
239.	DT157Pilot	1.2204	.	.	.	15.00
240.	DT157Pilot	0.7068	.	.	.	2.60
241.	DT157Pilot	2.5560	.	.	.	40.00
242.	DT157Pilot	0.1512	.	6.00	.	.
243.	DT157Pilot	2.2884	100.00	20.00	.	.
244.	DT157Pilot	0.4158	.	.	.	2.60
245.	DT157Pilot	0.2928	.	20.00	.	.
246.	DT157Pilot	1.7202	.	.	.	1.60
247.	DT157Pilot	0.3294	.	.	63.00	.
248.	DT157Pilot	1.9344	.	.	58.00	.
249.	DT157Pilot	43.7760	.	.	.	100.00
250.	DT157Pilot	0.9558	.	.	.	40.00
251.	DT157Pilot	1.5210	.	4.00	35.00	6.00
252.	DT157Pilot	0.6426	.	.	65.00	.
253.	DT157Pilot	0.7788	.	.	.	35.00
254.	DT157Pilot	4.4340	29.00	.	.	0.60
255.	DT157Pilot	15.6582	100.00	.	18.00	8.00
256.	DT157Pilot	6.3480	.	.	.	2.50
257.	DT157Pilot	1.6632	.	.	.	80.00
258.	DT157Pilot	7.6152	.	5.00	100.00	.
259.	DT157Pilot	3.2760	.	.	.	20.00
260.	DT157Pilot	1.3860	.	.	100.00	.
261.	DT157Pilot	4.8633	.	32.00	.	.
262.	DT157Pilot	5.7782	.	.	62.00	4.00
263.	DT157Pilot	0.0037	.	.	0.40	.
264.	DT157Pilot	1.4552	.	2.00	26.00	.
265.	DT157Pilot	0.6030	.	.	.	0.70
266.	DT157Pilot	0.8245	.	6.00	.	1.40
267.	DT157Pilot	0.0991	.	1.00	10.00	.
268.	DT157Pilot	0.0730	.	.	.	0.13
269.	DT157Pilot	0.7068	.	.	88.00	.
270.	DT157Pilot	1.3986	100.00	.	.	.
271.	DT157Pilot	9.5040	.	8.00	100.00	.
272.	DT157Pilot	0.2325	.	.	20.00	.
273.	DT157Pilot	9.5397	100.00	.	.	2.00
274.	DT157Pilot	3.2640	.	41.00	48.00	.
275.	DT157Pilot	1.1678	.	0.60	.	6.00
276.	DT157Pilot	0.4800	.	25.00	.	.
277.	DT157Pilot	1.8600	.	.	.	2.00
278.	DT157Pilot	1.4997	22.00	5.00	.	.
279.	DT157Pilot	4.8165	.	.	100.00	7.00

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
280.	DT157Pilot	0.8160	.	23.00	.	.
281.	DT157Pilot	0.7466	.	20.00	.	.
282.	DT157Pilot	9.9900	.	.	90.00	.
283.	DT157Pilot	0.5940	.	.	.	5.00
284.	DT157Pilot	0.5346	.	32.00	.	.
285.	DT157Pilot	0.8190	.	.	.	13.00
286.	DT157Pilot	2.7444	.	.	100.00	.
287.	DT157Pilot	0.6696	.	.	84.00	.
288.	DT157Pilot	3.2484	.	9.00	.	10.00
289.	DT157Pilot	2.5164	.	20.00	.	18.00
290.	DT157Pilot	1.3392	.	.	100.00	.
291.	DT157Pilot	1.0860	.	.	100.00	0.20
292.	DT157Pilot	0.3206	.	.	.	4.40
293.	DT157Pilot	41.7180	.	.	100.00	.
294.	DT157Pilot	19.4700	.	.	.	16.00
295.	DT157Pilot	4.1681	.	10.00	.	3.00
296.	DT157Pilot	0.9130	.	.	.	1.00
297.	DT157Pilot	2.0910	.	.	72.00	3.00
298.	DT157Pilot	2.8928	.	.	.	14.00
299.	DT157Pilot	3.8766	30.00	.	100.00	.
300.	DT157Pilot	1.1046	100.00	.	35.00	.
301.	DT157Pilot	1.7305	.	.	.	100.00
302.	DT157Pilot	2.5212	100.00	.	100.00	.
303.	DT157Pilot	6.0848	.	.	.	50.00
304.	DT157Pilot	9.1800	95.00	.	.	.
305.	DT157Pilot	1.4310	.	.	.	6.00
306.	DT157Pilot	1.6083	.	.	40.00	.
307.	DT157Pilot	0.2317	.	.	.	5.00
308.	DT157Pilot	2.1249	10.00	.	.	50.00
309.	DT157Pilot	2.4960	.	.	100.00	.
310.	Natl_CARB_GTI	20.4000	.	.	.	.
311.	Natl_CARB_GTI	14.4000	.	.	.	.
312.	Natl_CARB_GTI	13.9850	.	.	.	.
313.	Natl_CARB_GTI	13.8000	.	.	.	.
314.	Natl_CARB_GTI	13.2000	.	.	.	.
315.	Natl_CARB_GTI	7.2000	.	.	.	.
316.	Natl_CARB_GTI	6.9000	.	.	.	.
317.	Natl_CARB_GTI	6.4950	.	.	.	.
318.	Natl_CARB_GTI	6.4750	.	.	.	.
319.	Natl_CARB_GTI	5.7000	.	.	.	.
320.	Natl_CARB_GTI	5.4000	.	.	.	.
321.	Natl_CARB_GTI	5.0000	.	.	.	.
322.	Natl_CARB_GTI	5.0000	.	.	.	.
323.	Natl_CARB_GTI	5.0000	.	.	.	.
324.	Natl_CARB_GTI	4.0000	.	.	.	.
325.	Natl_CARB_GTI	3.9000	.	.	.	.
326.	Natl_CARB_GTI	3.2000	.	.	.	.
327.	Natl_CARB_GTI	2.4000	.	.	.	.
328.	Natl_CARB_GTI	2.3860	.	.	.	.
329.	Natl_CARB_GTI	2.1540	.	.	.	.
330.	Natl_CARB_GTI	2.1000	.	.	.	.
331.	Natl_CARB_GTI	2.0970	.	.	.	.
332.	Natl_CARB_GTI	2.0000	.	.	.	.
333.	Natl_CARB_GTI	1.9440	.	.	.	.
334.	Natl_CARB_GTI	1.8000	.	.	.	.
335.	Natl_CARB_GTI	1.5340	.	.	.	.
336.	Natl_CARB_GTI	1.5320	.	.	.	.
337.	Natl_CARB_GTI	1.4360	.	.	.	.
338.	Natl_CARB_GTI	1.3310	.	.	.	.
339.	Natl_CARB_GTI	1.2840	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
340.	Natl_CARB_GTI	1.2000	.	.	.	.
341.	Natl_CARB_GTI	1.2000	.	.	.	.
342.	Natl_CARB_GTI	1.1680	.	.	.	.
343.	Natl_CARB_GTI	1.0670	.	.	.	.
344.	Natl_CARB_GTI	1.0000	.	.	.	.
345.	Natl_CARB_GTI	1.0000	.	.	.	.
346.	Natl_CARB_GTI	0.8840	.	.	.	.
347.	Natl_CARB_GTI	0.8530	.	.	.	.
348.	Natl_CARB_GTI	0.8000	.	.	.	.
349.	Natl_CARB_GTI	0.7640	.	.	.	.
350.	Natl_CARB_GTI	0.7190	.	.	.	.
351.	Natl_CARB_GTI	0.6540	.	.	.	.
352.	Natl_CARB_GTI	0.6510	.	.	.	.
353.	Natl_CARB_GTI	0.6170	.	.	.	.
354.	Natl_CARB_GTI	0.6130	.	.	.	.
355.	Natl_CARB_GTI	0.6010	.	.	.	.
356.	Natl_CARB_GTI	0.6000	.	.	.	.
357.	Natl_CARB_GTI	0.6000	.	.	.	.
358.	Natl_CARB_GTI	0.6000	.	.	.	.
359.	Natl_CARB_GTI	0.6000	.	.	.	.
360.	Natl_CARB_GTI	0.6000	.	.	.	.
361.	Natl_CARB_GTI	0.6000	.	.	.	.
362.	Natl_CARB_GTI	0.6000	.	.	.	.
363.	Natl_CARB_GTI	0.6000	.	.	.	.
364.	Natl_CARB_GTI	0.6000	.	.	.	.
365.	Natl_CARB_GTI	0.5850	.	.	.	.
366.	Natl_CARB_GTI	0.5150	.	.	.	.
367.	Natl_CARB_GTI	0.4670	.	.	.	.
368.	Natl_CARB_GTI	0.4620	.	.	.	.
369.	Natl_CARB_GTI	0.4520	.	.	.	.
370.	Natl_CARB_GTI	0.4350	.	.	.	.
371.	Natl_CARB_GTI	0.3780	.	.	.	.
372.	Natl_CARB_GTI	0.2930	.	.	.	.
373.	Natl_CARB_GTI	0.2760	.	.	.	.
374.	Natl_CARB_GTI	0.2550	.	.	.	.
375.	Natl_CARB_GTI	0.1740	.	.	.	.
376.	Natl_CARB_GTI	0.1720	.	.	.	.
377.	Natl_CARB_GTI	0.1660	.	.	.	.
378.	Natl_CARB_GTI	0.1560	.	.	.	.
379.	Natl_CARB_GTI	0.1480	.	.	.	.
380.	Natl_CARB_GTI	0.1200	.	.	.	.
381.	Natl_CARB_GTI	0.0780	.	.	.	.
382.	Natl_CARB_GTI	0.0630	.	.	.	.
383.	Natl_CARB_GTI	0.0550	.	.	.	.
384.	Natl_CARB_GTI	0.0410	.	.	.	.
385.	Natl_CARB_GTI	0.0070	.	.	.	.
386.	Natl_OTD_GTI	95.4000	.	.	.	.
387.	Natl_OTD_GTI	78.6000	.	.	.	.
388.	Natl_OTD_GTI	24.3000	.	.	.	.
389.	Natl_OTD_GTI	16.2000	.	.	.	.
390.	Natl_OTD_GTI	14.4000	.	.	.	.
391.	Natl_OTD_GTI	13.9852	.	.	.	.
392.	Natl_OTD_GTI	13.8000	.	.	.	.
393.	Natl_OTD_GTI	13.2000	.	.	.	.
394.	Natl_OTD_GTI	11.8800	.	.	.	.
395.	Natl_OTD_GTI	7.8000	.	.	.	.
396.	Natl_OTD_GTI	7.8000	.	.	.	.
397.	Natl_OTD_GTI	7.2000	.	.	.	.
398.	Natl_OTD_GTI	3.6000	.	.	.	.
399.	Natl_OTD_GTI	3.6000	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
400.	Natl_OTD_GTI	3.6000	.	.	.	.
401.	Natl_OTD_GTI	2.3864	.	.	.	.
402.	Natl_OTD_GTI	2.2500	.	.	.	.
403.	Natl_OTD_GTI	1.8000	.	.	.	.
404.	Natl_OTD_GTI	1.8000	.	.	.	.
405.	Natl_OTD_GTI	1.8000	.	.	.	.
406.	Natl_OTD_GTI	1.5317	.	.	.	.
407.	Natl_OTD_GTI	1.4000	.	.	.	.
408.	Natl_OTD_GTI	1.3310	.	.	.	.
409.	Natl_OTD_GTI	1.3307	.	.	.	.
410.	Natl_OTD_GTI	1.2000	.	.	.	.
411.	Natl_OTD_GTI	1.2000	.	.	.	.
412.	Natl_OTD_GTI	1.2000	.	.	.	.
413.	Natl_OTD_GTI	1.2000	.	.	.	.
414.	Natl_OTD_GTI	0.8526	.	.	.	.
415.	Natl_OTD_GTI	0.8000	.	.	.	.
416.	Natl_OTD_GTI	0.6000	.	.	.	.
417.	Natl_OTD_GTI	0.6000	.	.	.	.
418.	Natl_OTD_GTI	0.6000	.	.	.	.
419.	Natl_OTD_GTI	0.6000	.	.	.	.
420.	Natl_OTD_GTI	0.6000	.	.	.	.
421.	Natl_OTD_GTI	0.6000	.	.	.	.
422.	Natl_OTD_GTI	0.6000	.	.	.	.
423.	Natl_OTD_GTI	0.6000	.	.	.	.
424.	Natl_OTD_GTI	0.6000	.	.	.	.
425.	Natl_OTD_GTI	0.6000	.	.	.	.
426.	Natl_OTD_GTI	0.6000	.	.	.	.
427.	Natl_OTD_GTI	0.6000	.	.	.	.
428.	Natl_OTD_GTI	0.6000	.	.	.	.
429.	Natl_OTD_GTI	0.6000	.	.	.	.
430.	Natl_OTD_GTI	0.6000	.	.	.	.
431.	Natl_OTD_GTI	0.6000	.	.	.	.
432.	Natl_OTD_GTI	0.6000	.	.	.	.
433.	Natl_OTD_GTI	0.6000	.	.	.	.
434.	Natl_OTD_GTI	0.6000	.	.	.	.
435.	Natl_OTD_GTI	0.6000	.	.	.	.
436.	Natl_OTD_GTI	0.6000	.	.	.	.
437.	Natl_OTD_GTI	0.6000	.	.	.	.
438.	Natl_OTD_GTI	0.6000	.	.	.	.
439.	Natl_OTD_GTI	0.6000	.	.	.	.
440.	Natl_OTD_GTI	0.6000	.	.	.	.
441.	Natl_OTD_GTI	0.6000	.	.	.	.
442.	Natl_OTD_GTI	0.6000	.	.	.	.
443.	Natl_OTD_GTI	0.6000	.	.	.	.
444.	Natl_OTD_GTI	0.6000	.	.	.	.
445.	Natl_OTD_GTI	0.6000	.	.	.	.
446.	Natl_OTD_GTI	0.6000	.	.	.	.
447.	Natl_OTD_GTI	0.0439	.	.	.	.
448.	Natl_WSU_EDF	109.4722	.	.	.	.
449.	Natl_WSU_EDF	69.7000	.	.	.	.
450.	Natl_WSU_EDF	13.2912	.	.	.	.
451.	Natl_WSU_EDF	10.3518	.	.	.	.
452.	Natl_WSU_EDF	8.9751	.	.	.	.
453.	Natl_WSU_EDF	7.8412	.	.	.	.
454.	Natl_WSU_EDF	6.8595	.	.	.	.
455.	Natl_WSU_EDF	6.1869	.	.	.	.
456.	Natl_WSU_EDF	5.9035	.	.	.	.
457.	Natl_WSU_EDF	5.6127	.	.	.	.
458.	Natl_WSU_EDF	4.7525	.	.	.	.
459.	Natl_WSU_EDF	4.5909	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
460.	Natl_WSU_EDF	4.5446	.	.	.	.
461.	Natl_WSU_EDF	4.0154	.	.	.	.
462.	Natl_WSU_EDF	3.7319	.	.	.	.
463.	Natl_WSU_EDF	3.2224	.	.	.	.
464.	Natl_WSU_EDF	2.9655	.	.	.	.
465.	Natl_WSU_EDF	2.8829	.	.	.	.
466.	Natl_WSU_EDF	2.7585	.	.	.	.
467.	Natl_WSU_EDF	2.7051	.	.	.	.
468.	Natl_WSU_EDF	2.6703	.	.	.	.
469.	Natl_WSU_EDF	2.5612	.	.	.	.
470.	Natl_WSU_EDF	2.5575	.	.	.	.
471.	Natl_WSU_EDF	2.3693	.	.	.	.
472.	Natl_WSU_EDF	2.1575	.	.	.	.
473.	Natl_WSU_EDF	2.1111	.	.	.	.
474.	Natl_WSU_EDF	2.0437	.	.	.	.
475.	Natl_WSU_EDF	2.0330	.	.	.	.
476.	Natl_WSU_EDF	1.8371	.	.	.	.
477.	Natl_WSU_EDF	1.7565	.	.	.	.
478.	Natl_WSU_EDF	1.6219	.	.	.	.
479.	Natl_WSU_EDF	1.4635	.	.	.	.
480.	Natl_WSU_EDF	1.4562	.	.	.	.
481.	Natl_WSU_EDF	1.4498	.	.	.	.
482.	Natl_WSU_EDF	1.4265	.	.	.	.
483.	Natl_WSU_EDF	1.4140	.	.	.	.
484.	Natl_WSU_EDF	1.1119	.	.	.	.
485.	Natl_WSU_EDF	1.0910	.	.	.	.
486.	Natl_WSU_EDF	1.0518	.	.	.	.
487.	Natl_WSU_EDF	1.0454	.	.	.	.
488.	Natl_WSU_EDF	1.0454	.	.	.	.
489.	Natl_WSU_EDF	1.0444	.	.	.	.
490.	Natl_WSU_EDF	0.9838	.	.	.	.
491.	Natl_WSU_EDF	0.9539	.	.	.	.
492.	Natl_WSU_EDF	0.9414	.	.	.	.
493.	Natl_WSU_EDF	0.9081	.	.	.	.
494.	Natl_WSU_EDF	0.8880	.	.	.	.
495.	Natl_WSU_EDF	0.8552	.	.	.	.
496.	Natl_WSU_EDF	0.8482	.	.	.	.
497.	Natl_WSU_EDF	0.8213	.	.	.	.
498.	Natl_WSU_EDF	0.7985	.	.	.	.
499.	Natl_WSU_EDF	0.7646	.	.	.	.
500.	Natl_WSU_EDF	0.7400	.	.	.	.
501.	Natl_WSU_EDF	0.6807	.	.	.	.
502.	Natl_WSU_EDF	0.6678	.	.	.	.
503.	Natl_WSU_EDF	0.6671	.	.	.	.
504.	Natl_WSU_EDF	0.6594	.	.	.	.
505.	Natl_WSU_EDF	0.6271	.	.	.	.
506.	Natl_WSU_EDF	0.5978	.	.	.	.
507.	Natl_WSU_EDF	0.5907	.	.	.	.
508.	Natl_WSU_EDF	0.5627	.	.	.	.
509.	Natl_WSU_EDF	0.5143	.	.	.	.
510.	Natl_WSU_EDF	0.5140	.	.	.	.
511.	Natl_WSU_EDF	0.5082	.	.	.	.
512.	Natl_WSU_EDF	0.4997	.	.	.	.
513.	Natl_WSU_EDF	0.4348	.	.	.	.
514.	Natl_WSU_EDF	0.4168	.	.	.	.
515.	Natl_WSU_EDF	0.3936	.	.	.	.
516.	Natl_WSU_EDF	0.3855	.	.	.	.
517.	Natl_WSU_EDF	0.3819	.	.	.	.
518.	Natl_WSU_EDF	0.3796	.	.	.	.
519.	Natl_WSU_EDF	0.3765	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
520.	Natl_WSU_EDF	0.3763	.	.	.	.
521.	Natl_WSU_EDF	0.3762	.	.	.	.
522.	Natl_WSU_EDF	0.3631	.	.	.	.
523.	Natl_WSU_EDF	0.3582	.	.	.	.
524.	Natl_WSU_EDF	0.3447	.	.	.	.
525.	Natl_WSU_EDF	0.3445	.	.	.	.
526.	Natl_WSU_EDF	0.3392	.	.	.	.
527.	Natl_WSU_EDF	0.3311	.	.	.	.
528.	Natl_WSU_EDF	0.3188	.	.	.	.
529.	Natl_WSU_EDF	0.3097	.	.	.	.
530.	Natl_WSU_EDF	0.3056	.	.	.	.
531.	Natl_WSU_EDF	0.2970	.	.	.	.
532.	Natl_WSU_EDF	0.2948	.	.	.	.
533.	Natl_WSU_EDF	0.2916	.	.	.	.
534.	Natl_WSU_EDF	0.2916	.	.	.	.
535.	Natl_WSU_EDF	0.2847	.	.	.	.
536.	Natl_WSU_EDF	0.2847	.	.	.	.
537.	Natl_WSU_EDF	0.2826	.	.	.	.
538.	Natl_WSU_EDF	0.2736	.	.	.	.
539.	Natl_WSU_EDF	0.2632	.	.	.	.
540.	Natl_WSU_EDF	0.2619	.	.	.	.
541.	Natl_WSU_EDF	0.2602	.	.	.	.
542.	Natl_WSU_EDF	0.2552	.	.	.	.
543.	Natl_WSU_EDF	0.2494	.	.	.	.
544.	Natl_WSU_EDF	0.2429	.	.	.	.
545.	Natl_WSU_EDF	0.2428	.	.	.	.
546.	Natl_WSU_EDF	0.2424	.	.	.	.
547.	Natl_WSU_EDF	0.2290	.	.	.	.
548.	Natl_WSU_EDF	0.2263	.	.	.	.
549.	Natl_WSU_EDF	0.2234	.	.	.	.
550.	Natl_WSU_EDF	0.2233	.	.	.	.
551.	Natl_WSU_EDF	0.2196	.	.	.	.
552.	Natl_WSU_EDF	0.2147	.	.	.	.
553.	Natl_WSU_EDF	0.2112	.	.	.	.
554.	Natl_WSU_EDF	0.2108	.	.	.	.
555.	Natl_WSU_EDF	0.2059	.	.	.	.
556.	Natl_WSU_EDF	0.1987	.	.	.	.
557.	Natl_WSU_EDF	0.1945	.	.	.	.
558.	Natl_WSU_EDF	0.1934	.	.	.	.
559.	Natl_WSU_EDF	0.1905	.	.	.	.
560.	Natl_WSU_EDF	0.1892	.	.	.	.
561.	Natl_WSU_EDF	0.1890	.	.	.	.
562.	Natl_WSU_EDF	0.1854	.	.	.	.
563.	Natl_WSU_EDF	0.1848	.	.	.	.
564.	Natl_WSU_EDF	0.1822	.	.	.	.
565.	Natl_WSU_EDF	0.1812	.	.	.	.
566.	Natl_WSU_EDF	0.1808	.	.	.	.
567.	Natl_WSU_EDF	0.1798	.	.	.	.
568.	Natl_WSU_EDF	0.1742	.	.	.	.
569.	Natl_WSU_EDF	0.1716	.	.	.	.
570.	Natl_WSU_EDF	0.1714	.	.	.	.
571.	Natl_WSU_EDF	0.1713	.	.	.	.
572.	Natl_WSU_EDF	0.1712	.	.	.	.
573.	Natl_WSU_EDF	0.1676	.	.	.	.
574.	Natl_WSU_EDF	0.1662	.	.	.	.
575.	Natl_WSU_EDF	0.1658	.	.	.	.
576.	Natl_WSU_EDF	0.1628	.	.	.	.
577.	Natl_WSU_EDF	0.1623	.	.	.	.
578.	Natl_WSU_EDF	0.1614	.	.	.	.
579.	Natl_WSU_EDF	0.1610	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
580.	Natl_WSU_EDF	0.1600	.	.	.	.
581.	Natl_WSU_EDF	0.1595	.	.	.	.
582.	Natl_WSU_EDF	0.1582	.	.	.	.
583.	Natl_WSU_EDF	0.1580	.	.	.	.
584.	Natl_WSU_EDF	0.1549	.	.	.	.
585.	Natl_WSU_EDF	0.1544	.	.	.	.
586.	Natl_WSU_EDF	0.1524	.	.	.	.
587.	Natl_WSU_EDF	0.1498	.	.	.	.
588.	Natl_WSU_EDF	0.1496	.	.	.	.
589.	Natl_WSU_EDF	0.1486	.	.	.	.
590.	Natl_WSU_EDF	0.1471	.	.	.	.
591.	Natl_WSU_EDF	0.1413	.	.	.	.
592.	Natl_WSU_EDF	0.1392	.	.	.	.
593.	Natl_WSU_EDF	0.1367	.	.	.	.
594.	Natl_WSU_EDF	0.1344	.	.	.	.
595.	Natl_WSU_EDF	0.1323	.	.	.	.
596.	Natl_WSU_EDF	0.1309	.	.	.	.
597.	Natl_WSU_EDF	0.1127	.	.	.	.
598.	Natl_WSU_EDF	0.0999	.	.	.	.
599.	Natl_WSU_EDF	0.0987	.	.	.	.
600.	Natl_WSU_EDF	0.0964	.	.	.	.
601.	Natl_WSU_EDF	0.0948	.	.	.	.
602.	Natl_WSU_EDF	0.0896	.	.	.	.
603.	Natl_WSU_EDF	0.0848	.	.	.	.
604.	Natl_WSU_EDF	0.0847	.	.	.	.
605.	Natl_WSU_EDF	0.0847	.	.	.	.
606.	Natl_WSU_EDF	0.0826	.	.	.	.
607.	Natl_WSU_EDF	0.0804	.	.	.	.
608.	Natl_WSU_EDF	0.0802	.	.	.	.
609.	Natl_WSU_EDF	0.0798	.	.	.	.
610.	Natl_WSU_EDF	0.0790	.	.	.	.
611.	Natl_WSU_EDF	0.0770	.	.	.	.
612.	Natl_WSU_EDF	0.0761	.	.	.	.
613.	Natl_WSU_EDF	0.0738	.	.	.	.
614.	Natl_WSU_EDF	0.0706	.	.	.	.
615.	Natl_WSU_EDF	0.0635	.	.	.	.
616.	Natl_WSU_EDF	0.0608	.	.	.	.
617.	Natl_WSU_EDF	0.0577	.	.	.	.
618.	Natl_WSU_EDF	0.0575	.	.	.	.
619.	Natl_WSU_EDF	0.0574	.	.	.	.
620.	Natl_WSU_EDF	0.0555	.	.	.	.
621.	Natl_WSU_EDF	0.0530	.	.	.	.
622.	Natl_WSU_EDF	0.0524	.	.	.	.
623.	Natl_WSU_EDF	0.0469	.	.	.	.
624.	Natl_WSU_EDF	0.0456	.	.	.	.
625.	Natl_WSU_EDF	0.0438	.	.	.	.
626.	Natl_WSU_EDF	0.0416	.	.	.	.
627.	Natl_WSU_EDF	0.0413	.	.	.	.
628.	Natl_WSU_EDF	0.0408	.	.	.	.
629.	Natl_WSU_EDF	0.0403	.	.	.	.
630.	Natl_WSU_EDF	0.0397	.	.	.	.
631.	Natl_WSU_EDF	0.0397	.	.	.	.
632.	Natl_WSU_EDF	0.0392	.	.	.	.
633.	Natl_WSU_EDF	0.0392	.	.	.	.
634.	Natl_WSU_EDF	0.0377	.	.	.	.
635.	Natl_WSU_EDF	0.0369	.	.	.	.
636.	Natl_WSU_EDF	0.0364	.	.	.	.
637.	Natl_WSU_EDF	0.0306	.	.	.	.
638.	Natl_WSU_EDF	0.0292	.	.	.	.
639.	Natl_WSU_EDF	0.0249	.	.	.	.

	<b>Study</b>	<b>Leak Rate (scfh)</b>	<b>Conc. BH (% gas)</b>	<b>Conc. CIP (% gas)</b>	<b>Conc. SSS (% gas)</b>	<b>Conc. US (% gas)</b>
640.	Natl_WSU_EDF	0.0249	.	.	.	.
641.	Natl_WSU_EDF	0.0249	.	.	.	.
642.	Natl_WSU_EDF	0.0240	.	.	.	.
643.	Natl_WSU_EDF	0.0236	.	.	.	.
644.	Natl_WSU_EDF	0.0229	.	.	.	.
645.	Natl_WSU_EDF	0.0219	.	.	.	.
646.	Natl_WSU_EDF	0.0197	.	.	.	.
647.	Natl_WSU_EDF	0.0177	.	.	.	.
648.	Natl_WSU_EDF	0.0155	.	.	.	.
649.	Natl_WSU_EDF	0.0144	.	.	.	.
650.	Natl_WSU_EDF	0.0137	.	.	.	.
651.	Natl_WSU_EDF	0.0110	.	.	.	.
652.	Natl_WSU_EDF	0.0096	.	.	.	.
653.	Natl_WSU_EDF	0.0084	.	.	.	.
654.	Natl_WSU_EDF	0.0082	.	.	.	.
655.	Natl_WSU_EDF	0.0069	.	.	.	.
656.	Natl_WSU_EDF	0.0054	.	.	.	.
657.	Natl_WSU_EDF	0.0051	.	.	.	.
658.	Natl_WSU_EDF	0.0036	.	.	.	.
659.	Natl_WSU_EDF	0.0029	.	.	.	.

# References

1. Lamb, B., Edburg, S., Ferrara, T., Howard, T., Harrison M., Kolb, C., Townsend-Small, A., Dyck, W., Possolo, A., and Whetstone, J., *Direct Measurements Show Decreasing Methane Emissions from Natural Gas Local Distribution Systems in the United States*. Environmental Science & Technology, 2015. **49**: p. 5161-5169.
2. Farrag, K., Wiley, K., *Improving Methane Emission Estimates for Natural Gas Distribution Companies, Phase II – PE Pipes*. 2013, Operations Technology Development, OTD: Des Plaines, IL.
3. Moore, C., Stuver, S., and Wiley, K., *Final Report - Classification of Methane Emissions from Industrial Meters, Vintage vs Modern Plastic Pipe, and Plastic-lined Steel and Cast-Iron Pipe*. 2019: United States.
4. Ersoy, D., Adamo, M., Wiley, K., *Quantifying Methane Emissions from Distribution Pipelines in California*. 2019, California Air Resources Board (CARB).
5. *Gas Leak Abatement OIR, in Methane Leak Proceeding (R.15-01-008)*. 2015, California Public Utility Commission.
6. *California Senate Bill No. 1371, in Senate Bill (SB) 1371 (Statutes 2014, Chapter 525) 2014*, State of California.
7. *California Air Resources Board (CARB)*. 2020; Available from: <https://ww2.arb.ca.gov>.
8. *Methane Emissions from the Natural Gas Industry - Volume 9: Underground Pipelines*. 1996, Gas Research Institute and Environmental Protection Agency, GRI-94/0257.25, EPA-600/R-96-080: United States.
9. *Petroleum and Natural Gas Systems, in Title 40 Chapter 1 Subchapter C Part 98 Subpart W*, U.S.E.P. Agency, Editor. 2020 (Current Revision): United States.
10. *GRI-GHGCalc v1.0; GRI-99/0086*. 1999, Gas Technology Institute: Des Plaines, IL.
11. *Rules Governing Design, Construction, Testing, Maintenance, and Operation of Gas Gathering, Transmission, and Distribution Piping Systems, in General Order No. 112-F*. 2015, California Public Utilities Commission.
12. *GPTC Guide for Transmission and Distribution Piping Systems: 2012 Edition, in Guide Material Appendix G-192-11 Section 4.4 and 5.5*. 2012, Gas Piping Technology Committee (GPTC) of AGA.
13. Meeker, W.Q., Hahn, G.J., and Escobar, L.A., *Statistical Intervals: A Guide for Practitioners and Researchers*. 2 ed. 2017, Hoboken, NJ: Wiley.
14. Kupper, L.L., and Hafner, K.B., *How appropriate are popular sample size formulas?* American Statistician, 1989. **43**: p. 101-105.

15. Bolstad, W., *Introduction to Bayesian Statistics*. 2 ed. 2007, Hoboken New Jersey: Wiley-Interscience.
16. Ferson, S., *Bayesian methods in risk assessment, for Bureau de Recherches Geologiques et Minières*. Applied Biomathematics, 2005.
17. Kruschke, J.K., *Doing bayesian data analysis : a tutorial with R and BUGS*. 2011, Burlington, MA: Academic Press.
18. Elfron, B., Tibshirani, R.J., *An introduction to the bootstrap*. 1994, New York: Chapman and Hall.
19. Good, P.I., *Resampling Methods - a Practical Guide to Data Analysis*. 3 ed. 2006, Berlin: Birkhauser Publishing.
20. Lunneborg, C.E., *Data Analysis by Resampling: Concept and Applications*. 2000, Duxbury California.
21. D'Agostino, R.B., Belanger, A.J., and D'Agostino, Jr., R.B., *A suggestion for using powerful and informative tests of normality*. American Statistician, 1990. **44**: p. 316-321.
22. Royston, P., *An extension of Shapiro and Wilks's W test for normality to large samples*. Applied Statistics, 1982. **31**: p. 115-124.
23. Kolmogorov, A.N., *Sulla determinazione empirica di una legge di distribuzione*. Giornale dell'Istituto Italiano degli Attuari, 1933. **4**: p. 83-91.
24. Conover, W.J., *Practical Nonparametric Statistic*. 3 ed. 1999, New Yourk: Wiley.
25. Geyer, C.J., *Handbook of Markov Chain Monte Carlo*. Introduction to Markov chain Monte Carlo. 2011, Boca Raton, FL: Chapman & Hall/CRC.
26. Gelman, A., Gilks, W.R., and Roberts, G.O., *Weak convergence and optimal scaling of random walk Metropolis algorithms*. Annals of Applied Probability, 1997. **7**: p. 110-120.
27. Haario, H., Saksman, E., and Tamminen, J., *An adaptive Metropolis algorithm*. Bernoulli, 2001. **7**: p. 223-242.
28. Gelfand, A.E., Hills, S.E., Racine-Poon, A., and Smith, A.F.M., *Illustration of Bayesian inference in normal data models using Gibbs sampling*. Journal of the American Statistical Association **85**: 972-985., 1990. **85**: p. 972-985.
29. Savage, L., *Foundations of Statistics*. 1954, New York: John Wiley & Sons.
30. Hendrick, M.F., Ackley, R., Sanaie-Movahed, B., Tang, X., and Phillips, N.G., *Fugitive methane emissions from leak-prone natural gas distribution infrastructure in urban environments*. Environmental Pollution, 2016. **213**: p. 710-716.
31. Howard, T., *University of Texas study underestimates national methane emissions at natural gas production sites due to instrument sensor failure*. Energy Science & Engineering, 2015. **3**(5): p. 443-455.

32. Alvarez, R.A., Lyon, D.R., Marchese, A.J., Robinson, A.L., and Hamburg, S.P., *Possible malfunction in widely used methane sampler deserves attention but poses limited implications for supply chain emission estimates*. Elem Sci Anth, 2016. **4**.
33. *Hi-Flow Sampler Instruction Manual*. 2019; Available from: <https://www.mybacharach.com/wp-content/uploads/2015/08/0055-9017-Rev-7.pdf>
34. Van Hauwermeiren M, V.D.a.V.B.S., *A Compendium of Distributions*. 2 ed. 2012, Ghent, Belgium: Vose Software.
35. *Log-normal distribution*. 2019; Available from: [https://en.m.wikipedia.org/wiki/Log-normal\\_distribution](https://en.m.wikipedia.org/wiki/Log-normal_distribution).
36. Vose, D., *Risk Analysis: A quantitative guide*. 3 ed. 2008, West Sussex, England: John Wiley and Sons.
37. Schwartz, E.S., *The Stochastic Behavior of Commodity Prices: Implications for Valuation and Hedging*. J Finance, 1997. **52(3)**: p. 923-973.
38. Akaike, *A new look at the statistical model identification*. IEEE Trans. Automat. Control. , 1974. **vAC-19**: p. 716-723.
39. Hannan, E.J.a.Q., G.G., *The determination of the order of an autoregression*. J.R. Statistic. Soc. B, 1979. **41**: p. 190-195.