Computing Geologically Consistent Models from Geophysical Data

by

 $Justin \ Granek$

B.Sc., Acadia Universiy, 2009

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF SCIENCE

 in

The Faculty of Graduate Studies

(Geophysics)

THE UNIVERSITY OF BRITISH COLUMBIA

(Vancouver)

November 2011

 \bigodot Justin Granek 2011

Abstract

In this thesis an attempt is made to develop a methodology by which the information provided by downhole physical property logs can be leveraged to assist in the creation of constraints for the inversion of surface geophysics. I first motivate the research with an introduction to the utility of downhole physical property logging, including an overview of the diverse methods and data which can be acquired. Background information is also provided on statistical classification techniques and the UBC-GIF (University of British Columbia Geophysical Inversion Facility) inversion formulation so that the methodology can be properly understood.

The introduced methodology differs from previous attempts at incorporation of *a priori* information since it applies statistical classification of in situ physical property measurements (as opposed to physical property values inferred from geology) as the basis for constraints. Statistical classification, combined with the iterative nature of the scheme, act to propagate the information from the downhole physical property logs through-out the model with minimum user input required. This automated approach reduces the potential for bias from unsupported constraints, while maximizing the integration of the classification results.

The methodology is explained, and then demonstrated on three simple illustrative models. The results from these demonstrations are compared against unconstrained inversion, and the strengths and shortcomings of the methodology are discussed.

Table of Contents

A	ostra	ct ii
Ta	ble o	of Contents
Li	st of	Tables ix
Li	st of	Figures
1	Intr	oduction $\dots \dots \dots$
	1.1	Research Motivation
	1.2	Objectives
	1.3	Thesis Structure
2	Bor	ehole Geophysics
	2.1	Introduction
	2.2	Downhole Geophysical Methods
		2.2.1 Electrical
		2.2.1.1 Electrical Resistivity
		2.2.1.2 Spontaneous-Potential

		2.2.1.3	Single Point Resistance	13
		2.2.1.4	Induced Polarization	15
		2.2.1.5	Fluid Resistivity	15
		2.2.1.6	Inductive Conductivity	16
		2.2.1.7	Magnetic Susceptibility	17
	2.2.2	Radiome	tric	18
		2.2.2.1	Natural Gamma	18
		2.2.2.2	Gamma-Gamma Density	19
		2.2.2.3	Neutron	20
	2.2.3	Seismic		21
		2.2.3.1	Full Waveform Sonic	21
	2.2.4	Structur	al	22
		2.2.4.1	3-Arm Caliper	22
		2.2.4.2	Orientation Probe	23
		2.2.4.3	Temperature	25
		2.2.4.4	Acoustic Televiewer	25
		2.2.4.5	Optical Televiewer	26
2.3	Practio	cal Consid	lerations	27
	2.3.1	Casing		27
		2.3.1.1	Un-Cased Hole	27
		2.3.1.2	Steel Cased Hole	28
		2.3.1.3	Plastic Cased Hole	28
	2.3.2	Cost Ana	alysis	29
2.4	Summ	ary		29

Table of Contents

3	Stat	istical	Classification	31
	3.1	Introd	uction	31
		3.1.1	Supervised Learning	31
		3.1.2	Unsupervised Learning	32
		3.1.3	Classification	32
		3.1.4	Numerical Example	33
	3.2	K-Mea	ans Algorithm	33
		3.2.1	Numerical Example: K-Means Algorithm	34
	3.3	Expect	tation Maximization Algorithm	36
		3.3.1	Gaussian Mixture Model	36
		3.3.2	Maximum-Likelihood	38
		3.3.3	E-Step	39
		3.3.4	M-Step	ł1
		3.3.5	Convergence	13
		3.3.6	Numerical Example: EM Algorithm	14
			3.3.6.1 E-Step	14
			3.3.6.2 M-Step	ł5
			3.3.6.3 Convergence	ł7
		3.3.7	Practical Issues	ł7
			3.3.7.1 Number of Clusters	1 7
			3.3.7.2 Initialization	<u>1</u> 9
			3.3.7.3 Regularization	1 9
	3.4	Geoph	ysical Application	60

v

4	Geo	physic	cal Inversion	52
	4.1	Introd	uction	52
		4.1.1	Forward Problem	52
		4.1.2	Inverse Problem	53
	4.2	Mathe	ematical Formulation	54
		4.2.1	Data Misfit	54
		4.2.2	Objective Function	55
		4.2.3	Solving	56
		4.2.4	Control Parameters	57
			4.2.4.1 Weighting Functions	58
			4.2.4.2 Reference Models	60
			4.2.4.3 Bounds	60
	4.3	DCIP	2D	60
5	Iter	ative I	Inversion Technique	63
	5.1	Metho	odology	63
		5.1.1	Creating Geological Models	64
		5.1.2	Creating Geophysical Models	65
		5.1.3	Blind Inversion	69
		5.1.4	Simulating Physical Property Logs	71
		5.1.5	Classifying Downhole Data	72
		5.1.6	Creating Constraints	74

		5.1.6.1	Generating Bounds	74
		5.1.6.2	Generating a Reference Model	76
	5.1.7	Inversion	n With Constraints	76
	5.1.8	Updatin	g Constraints, Re-inverting & Iterating	78
		5.1.8.1	Classifying the Inversion Model & Updating the Bounds	78
		5.1.8.2	Re-inverting & Iterating	84
5.2	Model	I: Cylind	ler in a Half-space	86
	5.2.1	The Mo	del	86
	5.2.2	Downho	le Physical Property Logs & Classification	87
	5.2.3	Creating	g & Applying Constraints	89
	5.2.4	Updatin	g Constraints and Final Models	91
5.3	Model	II: Vertie	cal Contact with Resistive Overburden	98
	5.3.1	The Mo	del	98
	5.3.2	DC Data	a & Blind Inversion	101
		5.3.2.1	Data	101
		5.3.2.2	Inversion	102
	5.3.3	Downho	le Physical Property Logs & Classification	104
	5.3.4	Creating	g & Applying Constraints	106
	5.3.5	Updatin	g Constraints and Final Models	109
5.4	Model	III: Vert	ical Contact with Conductive Overburden	116
	5.4.1	The Mo	del	116

		5.4.2	DC Data & Blind Inversion
			5.4.2.1 Data
			5.4.2.2 Inversion
		5.4.3	Downhole Physical Property Logs & Classification 121
		5.4.4	Creating & Applying Constraints
		5.4.5	Updating Constraints and Final Models
6	Dise	cussion	134
Ū	61	Summ	arv of Results 134
	0.1	Jumm	
		6.1.1	Model I: Cylinder in a Halfspace
		6.1.2	Model II: Vertical Contact with Resistive Overburden 138
		6.1.3	Model III: Vertical Contact with Conductive Overburden
	6.2	False (Classification
	6.3	Choos	ing Parameters
	6.4	Real I	Data
		6.4.1	Extra Considerations
		6.4.2	Sample Classification
7	Con	clusio	ns
Bi	bliog	graphy	

List of Tables

2.1	Quick Reference for Downhole Geophysical Methods	10
3.2	γ^1_{ij} Values	45
3.3	Parameter Values	46
3.4	Convergence of the EM Method	47
4.1	Tunable Parameters in UBC GIF 3D Inversion Code	58
5.1	Model I: Physical Property Values of Geological Units	66
5.2	Model II: Physical Property Values of Geological Units	99
5.3	Model III: Physical Property Values of Geological Units	116

2.1	Sample of downhole geophysical data. Log abbreviations are listed on top, with units listed on bottom	4
2.2	Typical downhole geophysics setup (taken from Killeen [1997])	5
2.3	Diagram illustrating various downhole probes (taken from Killeen [1997])	7
2.4	Logging time: core logging vs hole logging (taken from Killeen [1997])	8
2.5	Square AC waveform with labels indicating during which por- tion of the signal Resistivity (R), Single-Point Resistance (SPR), Induced Polarization (IP), and Spontaneous-Potential (SP) are collected.	11
2.6	Example of a downhole probe for measuring electrical properties (Mount Sopris Poly-Electric Probe)	12
2.7	Schematic of typical circuit used for Spontaneous Potential/Single Point Resistance probe	e 14
2.8	Diagram of typical Induction Probe	16
2.9	Typical scintillation detector	18
2.10	Schematic of Compton scattering & photoelectric effect	19
2.11	Diagram of elastic scattering	20
2.12	Typical configuration of a seismic probe. Pulse is emitted from Tx, and travels through the hole walls to receivers at Rx1 and Rx2.	21

х

2.13	Motion of S-waves & P-waves	22
2.14	Diagram depicting the potential deviation between the expected and the true borehole traces	24
2.15	Sample of image from a televiewer log, with associated inter- pretaion	26
3.1	Simple 1D data example	33
3.2	Simple 1D data example	35
3.3	Simple 1D data example	44
3.4	Toy problem: EM algorithm results	47
3.5	Plot of number of clusters vs. Akaike information criterion	48
4.1	Diagram illustrating forward modeling	53
4.2	Diagram illustrating inversion	53
4.3	Diagram of the Tikhonov curve used for solving geophysical inverse problem	57
4.4	Unknown 1D model which connects points A and B, with prior knowledge of the values at points x_1 and x_2	58
4.5	Weighting function, W_x used to emphasize known points in the model. Larger values will be more penalized for not re- sembling the reference model	59
4.6	Example of 1D reference model, m_{ref}	59
5.1	Methodology flow chart	64
5.2	Geological model with two units (red and blue) \ldots	65
5.3	Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density	66

5.4	True model	67
5.5	Pole-dipole survey. Top with pole current electrode (red) on left, bottom with pole current electrode (red) on right	68
5.6	DC resistivity data (top) and associated percent errors (bottom)	69
5.7	Blind inversion, with (bottom) and without (top) horizontal surface weighting	70
5.8	Physical property logs. Properties listed above, units listed below, for three holes (listed at top)	72
5.9	Akaike information criterion vs number of clusters	73
5.10	Classification results. Left to right: rock type vs depth & scatter plot	73
5.11	Initial resistivity bounds constraints: upper bounds (top) and lower bounds (bottom)	75
5.12	Reference resistivity model	76
5.13	Recovered model from inversion with borehole constraints $\ . \ .$	77
5.14	1D diagram illustrating the model classification procedure. a) Binning the recovered model into rock types based on sta- tistical classification results. b) Resulting upper and lower bounds. c) Expansion coefficients. d) Resulting upper and lower bounds when b) is multiplied by c).	79
5.15	Model classification procedure. From top to bottom: previous recovered model, upper resistivity bounds without expansion coefficients, expansion coefficients, upper resistivity bounds with expansion coefficients.	81
5.16	Updated upper (top) and lower (bottom) bounds, iteration 1	83
5.17	Recovered model, iteration 1	84
5.18	Recovered model, iteration 3 (top) compared to recovered model from blind inversion result (bottom)	85
		xii

5.19	Model I: True model	86
5.20	Model I: DC resistivity data (top) and errors (bottom) $\ . \ . \ .$	86
5.21	Model I: Blind Inversion with (bottom) and without (top) horizontal surface weighting	87
5.22	Model I: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)	88
5.23	Model I: Akaike information criterion vs number of clusters .	88
5.24	Model I: Classification Results. Rock type vs depth (left) & scatter plot of physical property values (right)	89
5.25	Model I: Initial constraints. From top to bottom: upper re- sistivity bounds, lower resistivity bounds, and reference model.	90
5.26	Model I: Recovered model from inversion with borehole con- straints	91
5.27	Model I: Updated bounds, iteration 1	92
5.28	Model I: Recovered model, iteration 1	93
5.29	Model I: Updated bounds, iteration 2	94
5.30	Model I: Recovered model, iteration 2	95
5.31	Model I: Updated bounds, iteration III	96
5.32	Model I: Recovered model, iteration III (top), compared to blind inversion result (bottom)	97
5.33	Models II & III: Geological model with three units (red, green and blue)	98
5.34	Model II: Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density	100
5.35	Model II: True model	101

5.36	Model II: DC resistivity data (top) and associated percent errors (bottom)	102
5.37	Model II: Blind Inversion with (bottom) and without (top) horizontal surface weighting	103
5.38	Model II: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)	104
5.39	Model II: Akaike information criterion vs number of clusters .	105
5.40	Model II: Classification results. Rock type vs depth (left) & scatter plot of physical properties (right)	105
5.41	Model II: Initial constraints. From top to bottom: upper resistivity bounds, lower resistivity bounds, and reference model.	107
5.42	Model II: Recovered model from blind inversion with full sur- face weighting	108
5.43	Model II: Recovered model from inversion with borehole con- straints and surface weighting	108
5.44	Model II: Updated bounds, iteration 1	110
5.45	Model II: Recovered model, iteration 1	111
5.46	Model II: Updated bounds, iteration 2	112
5.47	Model II: Recovered model, iteration 2	113
5.48	Model II: Updated bounds, iteration 3	114
5.49	Model II: Recovered model, Iteration 3 (top) compared to blind inversion result (bottom)	115
5.50	Model III: Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density	117
5.51	Model III: True model	118
5.52	Model III: DC resistivity data (top) and associated percent errors (bottom)	119

5.53	Model III: Recovered model from blind inversion with (bot- tom) and without (top) horizontal surface weighting 12	20
5.54	Model III: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)	21
5.55	Model III: Akaike information criterion vs number of classes . 12	22
5.56	Model III: Classification results. Rock type vs depth (left) & scatter plot (right)	22
5.57	Model III: Initial constraints. From top to bottom: upper resistivity bounds, lower resistivity bounds, and reference model.12	24
5.58	Model III: Recovered model form blind inversion with full sur- face weighting	25
5.59	Model III: Recovered model from inversion with borehole con- straints	26
5.60	Model III: Updated bounds, iteration 1	27
5.61	Model III: Recovered model, iteration 2	28
5.62	Model III: Updated bounds, iteration $2 \dots $	29
5.63	Model III: Updated bounds, iteration $3 \dots \dots \dots \dots \dots \dots 13$	30
5.64	Model III: Recovered model, iteration 2	31
5.65	Model III: Recovered model, iteration 3	31
5.66	Model III: Updated bounds, iteration 4	32
5.67	Model III: Recovered model, iteration 4 (top), compared to blind inversion result (bottom)	33
6.1	Model I: Results. From top to bottom: True model, model re- covered from blind inversion, model recovered from suggested methodology. Outline of true lithological boundaries visible in black	36

6.2	Model I: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5	137
6.3	Model II: Results. From top to bottom: True model, model re- covered from blind inversion, model recovered from suggested methodology. Outline of true lithological boundaries visible in black	139
6.4	Model II: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5	140
6.5	Model III: Results. From top to bottom: True model, model recovered form blind inversion, model recovered from sug- gested methodology. Outline of true lithological boundaries visible in black	142
6.6	Model III: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5	143
6.7	Sample of real physical property logs used in classification	147
6.8	Akaike Information Criterion vs. number of clusters for sam- ple of real downhole data	148
6.9	Results from EM classification of sample of real downhole data. Left, scatter plot of values for magnetic susceptibil- ity vs. P-wave slowness vs. 8" normal resistivity. Right, plot of rock type vs depth	148
6.10	Classification results for sample of real downhole data	149

Acknowledgments

First and foremost I want to thank Doug Oldenburg for his support, encouragement, and enthusiasm in guiding me forward in my research. His assistance in focusing my efforts and helping me to navigate the ever changing research landscape was immeasurable, and immensely appreciated \langle

Thanks also to Laurens Beran for his help and guidance in understanding some of the more complicated aspects of statistical classification. An additional mention must be made for the support provided by the members of the UBC-GIF group.

DGI Geoscience Inc generously provided assistance with the documentation of borehole logging techniques, as well as supplied data for reference in the creation of synthetic models.

Chapter 1

Introduction

1.1 Research Motivation

The difficulties in finding economically viable mineral deposits has motivated the development of new exploration methodologies. This has led to greater efforts from the geophysical community to incorporate available sources of geological and geophysical information. Since the suite of available data types is diverse, the synthesis of multiple sources of information into a single coherent model can present many difficulties. In particular, the incorporation of geological constraints in the inversion of geophysical data has been investigated by various researchers (LeLievre [2009], Oldenburg and Pratt [2007], P.K. Fullagar [2007, 2008], Williams [2008]).

While valuable information can be gleaned from geological data, a challenge remains due to the disconnect between geological units and geophysical property values. Though descriptive, a distinct geological unit is not always able to uniquely characterize the physical properties of a volume of earth, and vice versa. Interpretation and translation to and from geological and geophysical units can introduce bias based on the expert's experience.

Furthermore, as geophysical inversion moves to 3D and geometries become more complex, simple interpolation of sparse information runs the risk of imposing inconsistent or unsupported constraints. Conversely, methods which themselves are overly complex lose appeal since they can demand an extensive amount of information which is not always available.

Current methodologies which exist to incorporate geological and geophysical information into inversion typically suffer from at least one of the aforementioned difficulties: either they require the user to interpret physical property values from geological information, or else they require the user to define some range of influence for each measurement in the model. Imposing constraints on a model which are biased in one of these ways can lead to recovered models unsupported by the data.

1.2 Objectives

The aim of this thesis is to present a means of incorporating information from downhole physical property logs into the inversion of geophysical data which maximizes the application of data while minimizing the degree of required user input.

In doing so it is hoped that such a methodology will enhance existing inversion technology by providing a set of statistically-based constraints which will leverage the depth, accuracy and multi-dimensionality of physical property logs to increase the depth of investigation and resolution of surface geophysics.

1.3 Thesis Structure

This thesis draws from three main topics: borehole geophysics, statistical classification, and geophysical inversion. As such, a chapter has been devoted to each to provide the reader with sufficient background knowledge of the applied topics before presenting a new methodology.

Chapter 2 presents an introduction to borehole geophysics, beginning with the motivation for using borehole geophysics, including a basic explanation of how downhole geophysical logs are collected. This is followed by comprehensive¹ list of downhole geophysical methods, divided into sections based on the physical phenomenon. For each downhole method, basic relevant information is provided, such as what quantity is measured, how this is done, and a brief discussion of applications. Chapter 2 finishes with a short section discussing the merits and drawbacks of the various borehole casing methods, as well as a short comment on the cost of logging downhole physical properties.

¹The majority of common techniques are discussed

1.3. Thesis Structure

Chapter 3 introduces the mathematical background behind statistical classification used in this thesis. This is built up starting with a general introduction to statistical methods, followed by an explanation of the commonly used K-Means algorithm (Hastie et al. [2001]). A toy example is used to illustrate how classification occurs, which leads to the introduction of probabilistic modeling and the Expectation-Maximization algorithm (Dempster et al. [1977], Hastie et al. [2001], McLachlan and Krishnan [2008]). The mathematical framework of the EM algorithm is explained and demonstrated, using the same toy example from before. This chapter ends with a brief discussion of the practical issues of applying the EM algorithm to a data set, including choosing the number of clusters and initialization.

Chapter 4 gives a quick introduction to geophysical inversion, with focus placed on the basic mathematical framework of the methodology developed at the University of British Columbia Geophysical Inversion Facility, as per Oldenburg and Li [2005]. This is followed by a short discussion of the various control parameters the UBC-GIF inversion codes are equipped with, and their impact on the recovered model.

After providing this necessary background, Chapter 5 introduces a methodology for incorporating these three elements into a procedure whereby geophysical inversion of surface data can be constrained via classification of downhole physical property logs. This chapter begins with an overview of the methodology, followed by a detailed explanation of the process broken down into individual steps. The entire scheme is then demonstrated on three simple geological models by stepping through the procedure and presenting the results.

The thesis concludes with Chapters 6 & 7, in which a discussion of the results is presented, including a critical analysis of the merits and difficulties of the suggested methodology, and recommendations for further related research.

Chapter 2

Borehole Geophysics

2.1 Introduction

What is borehole geophysics?

Borehole geophysics is a branch of geophysics in which sensors are lowered into the earth via drill holes in order to obtain measurements of the in-situ physical rock properties (figure 2.1). Also referred to as downhole geophysics or well logging, the technique has many applications, ranging from natural resource exploration (ie: oil, gas or minerals) to geotechnical studies of the earth (Pickett [1970]).



Figure 2.1: Sample of downhole geophysical data. Log abbreviations are listed on top, with units listed on bottom

2.1. Introduction

How does borehole geophysics work?

A standard mineral exploration project will begin with exploration on the surface (including geology, geochemistry and geophysics). Once a target has been identified, strategic holes will be drilled in an attempt to intersect the ore body. The primary goal of the drill hole has traditionally been to recover the drill core, a sample column of rock extracted from the hole, on which analysis can be performed to determine which geological units were intersected.

Recently, it has begun to be more common for geoscientists to log drill holes as well. Once the hole has been drilled, it is possible for the geoscientists to lower probes into the well using a cable whose length is carefully measured in order to provide accurate measurements of depth (see figure 2.2).

Depending on the probe(s) being employed, different physical properties can be obtained. figure 2.3 below shows a schematic of the six main downhole configurations for data collection:



Figure 2.2: Typical downhole geophysics setup (taken from Killeen [1997])

Passive Source Measures energy emanating from the surrounding rocks.

- Active Source Supplies a source into the borehole to measure a response from the surrounding rocks.
- **Surface Source** Supplies a source at the surface to measure a response from the rocks surrounding the borehole.
- Fluid Properties Measures properties of the fluid within the borehole.
- $\label{eq:borehole walls by inspection using optical/acoustic televiewers.}$
- Mechanical Measurements Measure the location, size and condition of the hole itself.

In section 2.2 we will explore which physical properties can be measured using borehole geophysics, how they are collected, and some practical issues for the application of these techniques.



Figure 2.3: Diagram illustrating various downhole probes (taken from Killeen [1997])

7

Why use borehole geophysics?

Unlike conventional geophysical techniques, borehole geophysics allows for accurate measurements of physical rock properties at depth, creating a virtual window into the earth through which geoscientists can extract a wide range of different information.

In addition to this, many natural resource exploration programs will employ core logging as a standard component of their procedure. Thus more often than not the boreholes are drilled whether or not downhole geophysics are to be employed. Given that the majority of the cost in borehole geophysics results from the need to drill, downhole geophysics is a cost effective way to maximize the amount of data obtained from every dollar spent. Furthermore, the time spent to log the core is vastly greater than the time taken to log the hole, yet the data quality is nearly equal, as shown in figure 2.4:



Figure 2.4: Logging time: core logging vs hole logging (taken from Killeen [1997])

Having knowledge of physical rock properties can be very beneficial in constraining geophysical inversions of data from surface geophysics (Oldenburg and Pratt [2007]). How to best utilize this information is a topic of ongoing research that will be further explored in later chapters.

2.2 Downhole Geophysical Methods

This Section comprises an overview of a typical² suite of downhole geophysical methods, divided into categories based on the techniques or physics involved. For each method, the following are discussed:

- The basic quantity collected, and if possible, associated units.
- A brief overview of the underlying theory for data collection.
- Typical applications for this downhole geophysical survey.

It should be noted that this information is intended as an introduction to the methods, with basic explanations provided. For further information on all of the methods the reader is directed to the following resources: epa, gsc, usg.

Table 2.1 below is designed to serve as a summary of this section for quick reference, providing a list of the methods, the physical property they measure, common units used, and the some of the typical applications.

 $^{^{2}}$ It is recognized that many other methods exist. An attempt has been made in this section to introduce the most widely used methods, and to explain them in a general way.

${f Method}$	Physical Property	Units	Application
$\operatorname{Spontaneous}$	Resistivity	mV	Fluid flow
Potential			and salinity
Single Point	$\operatorname{Resistivity}$	Ω	Fractures
Resistance			and faults
$\operatorname{Inductive}$	$\operatorname{Conductivity}$	S/m	Map electrical
Conductivity			formations
Electrical	Resistivity	$\Omega \cdot m$	Map electrical
Resistivity			formations
Induced	Chargeability	mV/V	Disseminated
Polarization			ore bodies
Fluid	Resistivity	$\Omega \cdot m$	Porosity
Resistivity			and salinity
Magnetic	Magnetic	$SI \times 10^{-3}$	Magnetic
$\operatorname{Susceptibility}$	$\operatorname{Susceptibility}$		${ m minerals}$
Natural	Radioactivity	CPS	Lithology
Gamma			and alteration
Gamma-Gamma	Density	g/cm^3	Base metal
$\mathbf{Density}$			exploration
Neutron	Density	CPS	Porosity and
			moisture content
Full Waveform	Velocity	$\mu s/m$	Bulk elastic
Sonic			$\operatorname{properties}$
3-Arm	Diameter	mm	Hole diameter
Caliper			and condition
North-Seeking	Position	-	Georeferencing
Gyro			${ m measurements}$
Temperature	Temperature	°C	Fractures and
			fluid flow
Acoustic	-	-	Dip angles,
Televiewer			fractures and contacts
Optical	-	-	Dip angles,
Televiewer			fractures and contacts

2.2. Downhole Geophysical Methods

Table 2.1: Quick Reference for Downhole Geophysical Methods

2.2.1 Electrical

2.2.1.1 Electrical Resistivity

What is Collected? Voltage between two downhole electrodes of fixed spacing due to a dipole current source. Data are converted to resistivity, in $\Omega \cdot m$.

How is it Measured? Electrical Resistivity (R), or Normal Resistivity, is often measured in conjunction with the other electrical properties (Spontaneous-Potential, Single-Point Resistance, and Induced Polarization) using a square AC waveform (figure 2.5). The Resistivity data is collected during the constant current phase of the square wave on time.



Figure 2.5: Square AC waveform with labels indicating during which portion of the signal Resistivity (R), Single-Point Resistance (SPR), Induced Polarization (IP), and Spontaneous-Potential (SP) are collected.

Just like a surface dipole-dipole DC survey, resistivity is calculated by measuring the voltage between two potential electrodes (M & N) due to a dipole source from current electrodes (A & B). In figure 2.6 below, an electrical current, I is injected in the bottom electrode, labeled A, and travels to the mud plug, or surface electrode, B (not included in the figure). A potential is measured between electrodes M (M8, M16, M32 or M64 in the figure) and N (the cable armor - not labeled). Multiple M electrodes exist so as to provide different volumes of investigation, since this is related to the distance between the A and M electrodes (approximately a sphere of radius AM, centered around electrode A).

The apparent resistivity is calculated using Ohm's law

$$\frac{V}{I} = R = \frac{\rho \cdot l}{A} \tag{2.1}$$

where V is the voltage between electrodes MN, I is the current measured between electrodes AB, and R is the average resistance of the volume of investigation, which can be represented in terms of a resistivity, ρ , a current travel path, l, and a cross sectional area through which the current travels, A. Rearranging this equation, the resistivity can be represented as

$$\rho = G \cdot \frac{V}{I} \tag{2.2}$$

where the geometry is represented in terms of a single variable $G = \frac{A}{l}$, the geometric factor.



Figure 2.6: Example of a downhole probe for measuring electrical properties (Mount Sopris Poly-Electric Probe)

Applications: Resistivity is related to the porosity, salinity and metal content of the formation. Since the electrode spacing is known however, a quantitative assessment of resistivity is possible with a properly calibrated

instrument. Therefore the Resistivity log can be very useful in mapping conductors, faults, fractures, and determining bulk electrical properties of formation rocks. One drawback to be aware of however is that since the measurement is made between two potential electrodes, the signal can be affected deflection reversals caused by electrode spacing relative to bedding thickness.

2.2.1.2 Spontaneous-Potential

What is Collected? Potential difference between a surface electrode and a probe electrode, in mV.

How is it Measured? The Spontaneous-Potential (SP), or Self-Potential, is often measured in conjunction with the other electrical properties (Single Point Resistance, Resistivity, and Induced Polarization) using a square AC waveform (figure 2.5). The SP data is collected in the late off-time of the IP decay, after the injected current has fully dissipated. It is measured as the potential difference between the top probe electrode (M64 in figure2.6) and either a mud plug or the cable armor. The signal from the SP is related to the electrochemical potentials, electrokinetic potentials, and redox effects. Because of this, the measurement is highly dependent on the fluids within the borehole and thus can be highly variable depending on the environment.

Applications: Due to the relationship with fluid flow and salinity, SP is used mainly to determine lithology, bed thickness, formation fluid salinity, as well as permeability. It is important to be aware that SP measures the relative potential difference between the borehole fluid and the formation fluid, and therefore responses are very dependent on individual survey conditions such as borehole fluid salinity and hole diameter. Additionally, SP is particularly prone to high noise levels and anomalous deflections since it is sensitive to stray currents and equipment malfunction.

2.2.1.3 Single Point Resistance

What is Collected? Voltage between a surface electrode and a single probe electrode. Data are converted into resistance, in Ω .

How is it Measured? The Single-Point Resistance (SPR) is often measured in conjunction with the other electrical properties (Spontaneous-Potential, Resistivity, and Induced Polarization) using a square AC waveform (figure 2.5). The SPR data - like the Resistivity data - is collected during the constant current phase of the square wave on time. Current is injected at electrode A and travels through the ground to electrode B (a mud plug), and a voltage is measured between the same two electrodes (see figure 2.7). Since the cross sectional area A is proportional to the AB spacing, as the probe gets further from the surface, the ratio of $\frac{l}{A}$ approaches zero, and therefore most of the response is due to resistive bodies close to either the probe electrode or the mud plug. Since the mud plug is stationary, any anomalous signal can be attributed to changes in resistivity near the probe electrode. SPR has the advantage that unlike the normal resistivity measurements, it is not afflicted with deflection reversals due to bedding thickness, and thus has higher vertical resolution.



Figure 2.7: Schematic of typical circuit used for Spontaneous Potential/Single Point Resistance probe

Applications: As a measure of resistivity, SPR is related to the porosity, salinity, and metal content of the formation, however since the travel path between current and voltage electrodes is unknown, the relationship cannot be quantified. As such, it is useful in determining the location of fractures

and faults, determining relative salinity of formation pore fluids, and as a metric for the grain size of conductive minerals.

2.2.1.4 Induced Polarization

What is Collected? Chargeability in mV/V.

How is it Measured? Induced Polarization (IP) is often measured in conjunction with the other electrical properties (Spontaneous-Potential, Single-Point Resistance, and Electrical Resistivity) using a square AC waveform (figure 2.5). The IP signal is measured at various time windows during the decay of the current from the on-time to the off-time. The standard secondary voltage measurement is taken in the middle of the off-time. The chargeability is defined as the ratio of this value to the primary on-time voltage.

Applications: The IP effect can be used to locate large disseminated ore bodies, or to indicate the presence of cation rich clays. Additionally, some alteration processes, such as pyritization, can provide a strong IP signal. While correlated with conductivity and resistivity logs, the IP log can sometimes produce very different results.

2.2.1.5 Fluid Resistivity

What is Collected? Resistivity of the borehole fluids in $\Omega \cdot m$.

How is it Measured? Fluids are passed through the probe where measurements are shielded from outside sources. As fluids pass through the probe, a resistivity measurement is taken using a DC Wenner array.

Applications: The Fluid Resistivity measurement is often used to correlate the temperature differences with differences in borehole fluids. It is also important as it can be used with Archie's law to estimate the porosity of rock units, as well as the resistivity of the rock (as opposed to the combination of fluid and rock resistivity).

2.2.1.6 Inductive Conductivity

What is Collected? Conductivity of formation rocks in S/m.

How is it Measured? Induction logging operates on the principles of Maxwell's Equations, similar to surface electromagnetic methods. The induction probe has two sets of coils, one to transmit a primary magnetic field into the earth, and another to receive the secondary induced field from the surrounding rocks (figure 2.8). An AC current in the transmitter coils gives rise to the primary magnetic field according to the Biot-Savart law. When this alternating magnetic field impinges on conductive bodies, eddy currents are induced in the bodies which then in turn give rise to secondary magnetic fields. The alternating secondary field will itself induce a current in the receiver coils which has two components, in-phase and quadrature. These components roughly correspond to Magnetic Susceptibility and Inductive Conductivity, and therefore these measurements are often collected together.



Figure 2.8: Diagram of typical Induction Probe

Applications: The main advantage of Inductive Conductivity is that unlike the rest of the electrical methods, it can operate in holes that are not filled with conductive fluids (ie: air filled, oil filled etc), or in plastic cased holes since the magnetic field can easily penetrate the casing. The Inductive Conductivity log also provides superior vertical resolution to many other electrical methods, and can be used to map conductive bodies, faults, pore fluids and other bulk electrical properties of formations. Many induction logging sondes have been designed to be insensitive to near field effects such as changes in hole diameter or fluid salinity, and therefore the majority of the response comes from 15-100cm out from the hole.

2.2.1.7 Magnetic Susceptibility

What is Collected? Magnetic Susceptibility of formation rocks in mCGS or $SI \times 10^{-3}$.

How is it Measured? Magnetic Susceptibility (MS) data is collected via induction logging, which operates on the principles of Maxwell's Equations, similar to surface electromagnetic methods. The induction probe has two sets of coils, one to transmit a primary magnetic field into the earth, and another to receive the secondary induced field from the surrounding rocks (figure 2.8). An AC current in the transmitter coils gives rise to the primary magnetic field according to the Biot-Savart law. When this alternating magnetic field impinges on conductive bodies, eddy currents are induced in the bodies which then in turn give rise to secondary magnetic fields. The alternating secondary field will itself induce a current in the receiver coils which has two components, in-phase and quadrature. These components roughly correspond to Magnetic Susceptibility and Inductive Conductivity, and therefore these measurements are often collected together.

Applications: The Magnetic Susceptibility of a formation is directly related to the quantity of magnetic minerals (magnetite and pyrhotite) contained within. A MS survey can therefore be a quick way to determine the amount of ferromagnetic minerals in a formation. Since changes in magnetic properties are often associated with hydrothermal alteration, MS can also be useful in mapping alteration zones, since magnetic minerals such as magnetite are oxidized to non-magnetic minerals such as hematite.

2.2.2 Radiometric

2.2.2.1 Natural Gamma

What is Collected? Counts per second of gamma rays of a given energy band, in *CPS* or *API*.

How is it Measured? Natural Gamma probes are typically equipped with either a sodium iodide or a cesium iodide scintillation detector which, when hit with gamma rays, gives off light. When paired with a photomultiplier tube and electronics (figure 2.9), this signal can be identified by its characteristic energy band as having come from the decay of a given radio-element.



Figure 2.9: Typical scintillation detector

The most common radio-elements encountered in natural environments are ${}^{40}K$, which decays to stable ${}^{40}Ca$ and ${}^{40}Ar$, and ${}^{214}Bi$ and ${}^{208}Tl$, which are radioactive daughter products of stable ${}^{238}U$ and ${}^{232}Th$, respectively. Since there should exist equilibrium between parent and daughter elements, it is possible to estimate the concentration of ${}^{238}U$ and ${}^{232}Th$ given the counts from each daughter product.

The data from a Natural Gamma probe will consist of at least the Total Count (gamma radiation counts per second hitting the scintillation detector), and often also include the windowed counts for each of ${}^{40}K$, ${}^{238}U$ and ${}^{232}Th$, with energy windows centered at 1.46MeV, 1.76MeV, and 2.62MeV, respectively.

Applications: Because the three main elements involved (potassium, uranium and thorium) are found in different concentrations in different rock types, the Natural Gamma log can be very informative when it comes to distinguishing different lithologies, particularly different clays, as well as detecting alteration zones. Since gamma rays can penetrate most mediums, the Natural Gamma log can be used in holes that have been cased with plastic or steel, or that are filled with mud, fluid or air. This versatility has made it one of the more popular downhole methods.

2.2.2.2 Gamma-Gamma Density

What is Collected? Bulk density of surrounding rocks, in g/cm^3

How is it Measured? The Gamma-Gamma Density probe is essentially a Natural Gamma probe with the addition of a weak radioactive source, such as ${}^{60}Co$, at the nose of the probe. The concept is the same, with an extra twist. The scintillation detector now counts gamma rays that have been back-scattered from the surrounding rocks. The density is derived from the ratio of two energy windows, typically one for low energies (< 200 keV), and the other for large energies. The number of counts in each of these windows is dependent on two effects: Compton scattering, and the photoelectric effect.



Figure 2.10: Schematic of Compton scattering & photoelectric effect

If the density of the surrounding rocks increases, more gamma rays will be scattered by the Compton effect, and thus the counts in both energy windows will decrease. If the atomic mass of the sampled rocks increases, however, then the photoelectric effect will absorb a larger portion of the scattered energy, and the low energy window will receive lower counts, while the high energy window will remain unaffected.
The ratio of high energy counts to low energy counts can therefore be treated as an indicator of atomic mass. If properly calibrated, the instrument can give good estimates of bulk density.

Applications: Because the Gamma-Gamma Density probe is particularly sensitive to heavy elements, it is well suited to detecting the presence of base metals in sampled rock since most rock forming minerals are relatively light. It can also be used successfully for lithological mapping in minerals with differing quantities of heavy elements such as iron and magnesium, or for detecting changes in porosity, water content or compaction.

2.2.2.3 Neutron

What is Collected? Gamma ray counts per second, in CPS.

How is it Measured? There exist three main variations of detectors, however they all operate on the same main principles. The Neutron probe employs a high-energy neutron source such as Americium-Beryllium, which sends fast neutrons into the surrounding rocks. These neutrons interact with the nuclei of the rocks by elastically scattering (figure 2.11), and emitting energy as they do so. Since the optimal collision occurs between two bodies of the same size, hydrogen is the optimal target for these fast neutrons, emitting gamma rays and slow neutrons which can be detected by the probe.



Figure 2.11: Diagram of elastic scattering

Applications: Since the Neutron probe is sensitive to hydrogen content in the sample rocks, it is ideally suited for measuring porosity and moisture content. It has also been used in a similar manner to the Natural Gamma probe to map lithology.

2.2.3 Seismic

2.2.3.1 Full Waveform Sonic

What is Collected? Compressional (P) and shear (S) wave velocities, in m/s.

How is it Measured? Seismic probes typically operate with at least one transmitter and two receivers (figure 2.12). The transmitter will pulse a P-wave into the hole, which will travel through the hole fluids as well as the hole walls, to the two receivers, one near and one far.



Figure 2.12: Typical configuration of a seismic probe. Pulse is emitted from Tx, and travels through the hole walls to receivers at Rx1 and Rx2.

S-waves, or shear waves, oscillate perpendicular to the direction of propagation, and can only travel through solids, whereas P-waves, or pressure waves, oscillate parallel to the direction of propagation, and can travel through all mediums (figure 2.13). The complete acoustic signal is recorded at both receivers, and given the character of each signal and the difference in arrival times of characteristic modes, an estimate of S and P wave velocities can be calculated.



Figure 2.13: Motion of S-waves & P-waves

Applications: When combined with surface seismic data, or used in a hole-to-hole configuration, seismic methods can be used for tomography. When the S and P wave velocities are combined with a density measurement, it is possible to estimate bulk elastic properties of the materials, such as Young's Modulus etc.

Additionally, seismic methods are useful in determining the porosity or permeability of rocks, or the location of fractures or faults within a formation. The penetration depth of the seismic waves is dependent on the frequency used by the transmitter, and is typically of the order of tens of centimeters.

2.2.4 Structural

2.2.4.1 3-Arm Caliper

What is Collected? Diameter of the borehole, in inches or cm.

How is it Measured? The 3-Arm Caliper is exactly that: it has three mechanical arms which measure the diameter of the hole. This is done by linking all three arms to a linear potentiometer with constant reference voltage. The DC output voltage from the potentiometer is converted to a frequency, which is then corrected so that it is approximately linearly proportional to borehole diameter.

Applications: This log is primarily used to help in the correlation and calibration of other instruments, since the response from many probes is dependent on hole diameter. It can also help to locate areas of fracturing or caving in an uncased hole, and can be used to proactively prevent damage to equipment sent downhole.

2.2.4.2 Orientation Probe

What is Collected? 3D position of the downhole probe.

How is it Measured? The Orientation probe can serve as a reference to delineate where the borehole goes once it is below the surface. Drill holes will often deviate in both azimuth and dip due to the heterogeneous nature of the material they are boring, and it can be difficult to predict where a hole will end up given the collar location.



Figure 2.14: Diagram depicting the potential deviation between the expected and the true borehole traces

Two main methods exist for tracking the position of the borehole. Three-Component Flux-gate Magnetometers continuously monitor the probe's relationship to the earth's magnetic field in order to determine the dip and azimuth. This can give high resolution positional data which can be easily filtered to remove any small scale local magnetic anomalies, however since the probe relies on magnetic fields, it can have difficulties operating in highly magnetic environments or inside of a steel cased hole.

The most recent orientation probes avoid this problem by using high accuracy North-seeking gyro-compasses. A gyro-compass is essentially a gyroscope with the addition of a component which applies torque whenever the axis is not pointing North. Rather than using the earth's magnetic field, the gyrocompass uses the earth's spin, and the conservation of angular momentum to maintain its positional accuracy. Because of this, north seeking gyrocompasses are not susceptible to difficulties in magnetic environments, and can be used inside a cased hole.

Applications: Orientation probes are used primarily for the georeferencing of other measurements within the borehole. By giving the three di-

mensional position of the borehole, all other measurements can be correctly located within the trace of the hole.

2.2.4.3 Temperature

What is Collected? Temperature of borehole fluids, in ^{o}C .

How is it Measured? The temperature probe uses high sensitivity thermistor beads to measure changes in temperature of the borehole fluids and surrounding rocks. The thermistor sends a digital signal to the surface, which when calibrated using the correct constants can be used to compute the temperature.

Applications: Changes in thermal properties of rocks mainly indicate cracks or changes in thermal conductivity. Cracks allow for fluid flow which can produce characteristic signals in the temperature log, whereas conductive mineralization can be seen as an increase in temperature when the thermal background is relatively quiet. Often a Temperature Gradient probe is used to further increase the resolution and amplify the signal. Temperature log-ging is typically carried out on the down run so that the acquired data is unaffected by the thermal signature of the instruments.

2.2.4.4 Acoustic Televiewer

What is Collected? A 360^o ultrasound image of the borehole walls.

How is it Measured? An Acoustic Televiewer operates by recording a 360° ultrasound image of the borehole walls (figure 2.15). This is achieved by emitting a high frequency (~1.2MHz) energy wave from the probe head, and recording the amplitude of the wave after it has reflected off of the borehole wall. The reflectivity of the borehole wall depends mainly on the impedance (product of density and acoustic velocity) contrast between the inside and outside of the hole wall.



Figure 2.15: Sample of image from a televiewer log, with associated interpretaion

Applications: The ultrasound image can be related to the density of the rock in the borehole wall. Since the Acoustic Televiewer records continuously with such high resolution, the resulting pseudo-3D image of the density structure of the hole can be used for determining dip angles, locating fractures, or lithological contacts.

2.2.4.5 Optical Televiewer

What is Collected? A 360° image of the borehole walls.

How is it Measured? The Optical Televiewer is essentially a video camera equipped with specially designed optics to record a 360° image of the borehole (figure 2.15). The tip of the probe is mounted with a light ring to supply the necessary light, and the image is recorded using a high resolution, high sensitivity CCD camera.

Applications: Similarly to the Acoustic Televiewer, the Optical Televiewer is well suited for determining dip angles, locating thin beds and fractures, and for lithological interpretation.

2.3 Practical Considerations

2.3.1 Casing

When drilling a borehole, one of the major considerations is whether or not to case the hole. Depending on the geology and the goal of the drilling program, it is sometimes advisable to case the holes for various reasons, including

- to prevent caving in of unconsolidated sediments
- to prevent loss of borehole fluids
- to prevent contamination of target/deposit fluids
- to facilitate the travel of probes within the hole
- to mitigate differences in downhole pressure

Casing typically comes in two materials: steel (common in the petroleum industry) or plastics, such as PVC. There are arguments for and against both materials, as well as for not casing the hole at all. Depending on if/how the hole is cased, certain geophysical methods might no longer be possible. A summary of the limiting effect of borehole casing is provided in the final column of table 2.1.

2.3.1.1 Un-Cased Hole

An uncased borehole runs the risk of caving in. This risk increases in unconsolidated environments, or environments which are subjected to high pressure (such as at depth). Additionally, over time even the best holes can begin to collapse due to natural events such as seismic activity, groundwater flow, or other unforeseen events. For applications such as groundwater wells or hydrocarbon recovery, an uncased borehole is less than ideal since the fluids will easily disperse through the hole walls and thus make surface recovery nearly impossible.

On the other hand, for exploration purposes, an uncased hole has the advantage of providing an uninhibited means to sample physical properties at depth. Since certain geophysical methods require contact with the borehole walls (ie: resistivity, televiewers, magnetic susceptibility), uncased holes are idea for downhole geophysical logging, and thus it is common to log holes either *while* drilling, or simply before the hole is cased (if it is to be cased at all).

2.3.1.2 Steel Cased Hole

At the other end of the spectrum is the option to case the hole in steel. This is common practice for hydrocarbon exploration as it provides a solid, sealed conduit through which to extract oil and gas. In addition to this, it increases the lifespan of a borehole by significantly decreasing the risk of a collapse or cross contamination of fluids. By using different grades of steel, as well as different thicknesses and diameters of pipe, a hole can be cased down to great depths, thus isolating the hole from the surrounding rocks/fluids.

The downside to casing a hole in steel is the limiting effect it has on downhole geophysics. As was previously mentioned, some geophysical techniques required contact with the borehole walls in order to perform proper measurements, and by casing the hole in steel certain methods will no longer be possible. In particular, methods which rely on electric or electromagnetic phenomenon (resistivity, inductive conductivity, magnetic susceptibility etc) are drastically effected by encasing the hole in a solid conductive body, and thus cannot be reliably collected.

2.3.1.3 Plastic Cased Hole

Casing the hole with plastic such as PVC piping can be seen as a middle ground between steel casing and no casing. While PVC is not as rigid or reliable as steel, it also does not suffer as badly from the same electromagnetic limitations. Since plastic is not conductive, methods which do not require direct contact with the borehole wall (inductive conductivity, magnetic susceptibility) can still be used. Because of this, plastic casing is commonly used in mineral exploration and geotechnical studies.

2.3.2 Cost Analysis

The costs of a drilling program can be prohibitively expensive. Traditionally, drilling has been the final stage of exploration program: surface and airborne measurements would be collected to locate the approximate location of the target, and only then would drilling begin in an attempt to strike it. Previously, rather than log holes, the emphasis has been to log the core extracted from the holes, and to focus available resources on drilling more holes.

The cost of logging each hole is a small fraction of the considerable investment made to drill, while the amount of information provided by downhole logging is extremely valuable. Combined with other sources of information, downhole geophysics can be used to inform the location of future drilling, thus saving considerable investment by avoiding misplaced holes. Additionally, downhole logging has several advantages over core logging, such as higher vertical resolution, more spatially relatable, and faster logging times (see figure 2.4).

Despite the added cost of mobilization and demobilization of logging equipment, downhole geophysical logging is still a minor expense when compared to the cost of drilling the hole, and even more so should the hole be cased - casing can almost double the cost of establishing a borehole, depending on the depth and casing material. As such, logging boreholes should become a standard practice, given the high returns on such a relatively small investment.

2.4 Summary

As one can see from this overview of borehole geophysics, applied properly, downhole logging can be an important tool for mineral exploration. The versatility of the numerous techniques combined with the accuracy achieved in measuring in situ physical property values can provide extensive information about the subsurface.

The difficulty with applying this information to geophysical inversion lies in the complexity of the downhole logs. Interpreting multiple logs measuring different physical properties, even in the same hole, can be extremely difficult. This is further complicated by adding multiple holes (possible in a three dimensional configuration), and noting that the high vertical resolution achieved by downhole logs (on the order of tens of centimeters) is too great for most other exploration techniques (surface geophysics, geological mapping etc).

For these reasons, incorporating such a highly dimensional, high resolution data set into geophysical inversion presents difficulties. As a first step, the information from the various physical property logs and holes can be synthesized into coherent, well defined parameters using statistical classification.

Chapter 3

Statistical Classification

3.1 Introduction

In the previous chapter downhole physical property logs were presented as a viable tool for geophysical exploration. One of the advantages discussed was the wide array of different methods and techniques sensitive to a variety of physical properties, including conductivity, magnetic susceptibility, density, and acoustic velocity.

Due to this versatility, despite the one-dimensional nature of a borehole trace, the data acquired from logging a hole can be highly dimensional. Such datasets lend themselves nicely to statistical methods of discrimination to explain the underlying patterns in the data.

These methods can be divided into two main categories: supervised learning and unsupervised learning. The distinction between the two lies in the existence or lack of training data, \tilde{X} . Training data can be defined as being a subset of the data, X, for which the response (in classification problems this might be defined as the class), Y, is already known.

3.1.1 Supervised Learning

Given training data, \tilde{X} , supervised learning methods will apply this information to 'train' the algorithm how to assign an output for the rest of the data-set, for which the response is unknown. Mathematically this can be expressed as the following two step process:

1. Train the mapping function F using training data X:

$$\tilde{Y} = F\left(\tilde{X}\right)$$

2. Apply mapping function F to solve for unknown outputs Y, since

$$Y = F(X)$$

For downhole physical property logs, supervised learning methods could be applied if a library of physical properties and the associated lithological unit exists. In such a case the mapping function could be trained to recognize certain combinations of physical property values, and then assign a lithological unit defined as having those properties.

3.1.2 Unsupervised Learning

In a more realistic scenario, such a library might not exist. In the absence of training data, unsupervised learning methods can be used in which other means of determining outputs are employed. To make the problem more explicit, what is really needed is a method for determining the mapping function F without knowing the relationship between any of the inputs (X) and outputs (Y).

Since most of the quantities involved in this problem are unknown, it is helpful to re-cast the problem in terms of statistical quantities. As such, the problem becomes one of finding the parameters of a distribution F which explain the data X. This distribution can, in general, be any function, and once it is discovered, it can be applied to determine the outputs, Y.

One way of simplifying this problem is to divide a large, complex data-set into a number of smaller, simpler data-sets, and to attempt to explain each of them with (relatively) simple functions. Known as *clustering*, such algorithms will attempt to group similar data together such that the difference between each cluster is greater than the difference between data *within* each cluster.

3.1.3 Classification

Within both supervised and unsupervised learning, one goal can be to *classify* a data-set. This implies that a certain number of *classes* exist which divide the data according to some metric. In the formulation

 $Y = F\left(X\right)$

Y would become the solution such that y_i states which class x_i came from, using F as the classifier. Taking this further, in the context of clustering algorithms, each cluster can be defined as a class, and thus by clustering the data two goals are accomplished: the classifier is defined (as a set of clusters with defined distributions) and each datum is assigned to a class (by cluster membership).

In order to better illustrate how such methods of statistical classification can be applied to a real data-set, consider the following toy problem.

3.1.4 Numerical Example

In the following figure a simple one dimensional data-set is presented, in which values listed above the line, j = 1, 2..7, are the indices of the data values listed below the line $x_j = 1, 2, 3, 4, 6, 7, 8$.

j =	1	2	3	4	5	6	7
	-0	-0	-0-	-0		-0	_0
$\frac{x}{x_j} =$	1	2	3	4	6	7	8

Figure 3.1: Simple 1D data example

By inspection, it is simple to separate the distribution of data into clusters: clearly there are two, with one encapsulating data j = 1 - 4, and the other data 5 - 7. The goal now is to show how a computer can be taught to apply the same reasoning to determine the correct class for each datum (with possible classes being either cluster 1 or cluster 2).

3.2 K-Means Algorithm

In order to solve for the cluster assignments for all data x_j , j = 1..7, in figure 3.1, the K-Means algorithm will be introduced and applied. One of the simplest and most common clustering algorithms in use today, the K-Means algorithm operates by taking a user defined number of clusters, K, and iteratively completing the following two steps until convergence:

- 1. Assign every datum to the cluster with the closest mean value
- 2. Recalculate the mean value of each of the clusters

The algorithm begins by randomly assigning every datum to one of the K clusters. The mean value of each cluster is then calculated as

$$\mu_k = \frac{1}{N_k} \sum_{x_i \exists C_k} x_i \tag{3.1}$$

where N_k is the number of data in the k^{th} cluster, and $x_i \exists C_k$ denotes these data which were assigned to the k^{th} cluster. Given the current set of K mean values, each datum is reassigned to the cluster with the closest mean value, as defined by the Euclidean norm

Distance between
$$x_i$$
 and $\mu_k = d_{ik} = ||x_i - \mu_k||^2$

such that the intra-cluster variance

$$\sum_{k=1}^{K} N_k \sum_{x_i \exists C_k} d_{ik} \tag{3.2}$$

is minimized. These two steps are iterated until convergence - when the cluster assignments no longer change.

3.2.1 Numerical Example: K-Means Algorithm

Applying the K-Means algorithm to solve our toy problem, lets start by correctly guessing that there are two clusters. This simplifies the problem to solving for the mean value of each cluster, and the corresponding cluster assignments. The algorithm will be initialized by randomly assigning each datum to one of the two clusters, which will be referred to as C_1 and C_2 :

j =	1	2	3	4	5	6	7
C_1							
C_2		\checkmark	\checkmark				

The mean value of each cluster can then be calculated using Eq. 3.1:

$$\mu_1 = \frac{1}{3} \left(1 + 4 + 7 \right) = \frac{12}{3} = 4$$
$$\mu_2 = \frac{1}{4} \left(2 + 3 + 6 + 8 \right) = \frac{19}{4} = 4.75$$

It can be shown that - as one would expect - the cluster assignments which minimize Eq. 3.2 are:

j =	1	2	3	4	5	6	7
C_1	\checkmark		\checkmark				
C_2							

and the corresponding mean values are

$$\mu_1 = \frac{1}{4} \left(1 + 2 + 3 + 4 \right) = \frac{10}{4} = 2.50$$
$$\mu_2 = \frac{1}{3} \left(6 + 7 + 8 \right) = \frac{21}{3} = 7.00$$

with a variance of

$$\sigma_1^2 = \frac{1}{4} \sum_{x_i \exists C_1} \|x_i - 2.5\| = 1.25$$
$$\sigma_2^2 = \frac{1}{3} \sum_{x_i \exists C_2} \|x_i - 7\| = 1.00$$

Therefore the solution can be shown as



Figure 3.2: Simple 1D data example

where the red diamonds indicate the mean values and the pale red ellipses denote the variance. K-Means has successfully arrived at the correct solution to this toy problem, however one can imagine certain complications that would not be handled quite so nicely:

- 1. K-Means creates spherical clusters with shared variance for all dimensions; it therefore might not be ideal for modeling more complex processes.
- 2. The algorithm is very simplistic in its deterministic cluster assignment; what if a datum lies directly between two clusters? Which is the *closest* then?

In order to handle such difficulties, the more general Expectation-Maximization algorithm for clustering will be introduced and applied.

3.3 Expectation Maximization Algorithm

The EM algorithm, having been first formally presented in 1977 in a paper by Arthur Dempster, Nan Laird, and Donald Rubin, is now a well known and well studied algorithm. As such, many resources exist to assist in the understanding of the underlying theory, and it should be noted that all derivations in this chapter have been re-written based on a survey of these sources. To make the derivation more tangible, only relevant details are presented in this chapter. For further information, the reader is directed to the following resources: Chen and Gupta [2010], Dempster et al. [1977], Do and Batzoglou [2008], Fraley and Raftery [1998], Hastie et al. [2001], Mclachlan and Peel [2000], McLachlan and Krishnan [2008], Schneider [2001].

In this section a more general formulation of the toy problem will be presented to assist in the derivation and explanation of the Expectation-Maximization (EM) algorithm. To make for a simple example, the toy problem will still be solved using the EM algorithm, followed by a discussion of the benefits of such a generalized formulation.

3.3.1 Gaussian Mixture Model

Let the data X now be generalized to take on the form of a K-dimensional cloud of N data points:

$$X = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_K^1 \\ x_1^2 & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ x_1^N & \dots & \dots & x_K^N \end{bmatrix}$$
(3.3)

with each datum having an input from each dimension k so that the j^{th} datum is:

$$x_j = [x_{j1} \ x_{j2} \ \dots \ x_{jK}]$$

In these terms, the one dimensional toy problem has K = 1 and N = 7. To facilitate the solving of this complex problem, it will be assumed that each datum x_j can be modeled as having been independently sampled from a generative distribution, or likelihood function. Written $l(x_j|\theta)$, and read "the likelihood of observing datum x_j given the unknown set of defining parameters θ ", where θ is a set of parameters which we hope to discover from the observed data X. In general, $l(x_j|\theta)$ can be any distribution, however due to its inherent flexibility in fitting most functions, a Gaussian Mixture Model (GMM) is often used:

$$l(x_j|\theta) = \sum_{i=1}^{M} \alpha_i g_i(x_j|\theta_i) \quad \text{with} \quad \sum_{i=1}^{M} \alpha_i = 1$$
(3.4)

where M is the number of Gaussian components and $\theta_i = [\mu_i, \Sigma_i]$ is now taken to be the mean value and covariance of the i^{th} K-dimensional Gaussian distribution:

$$g_i(x_j|\theta_i) = \frac{1}{(2\pi)^{\frac{K}{2}} \det(\Sigma_i)^{\frac{1}{2}}} exp\left(\frac{-1}{2} (x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i)\right)$$
(3.5)

with α_i as the weight given to the i^{th} distribution. The constraint that

$$\sum_{i=1}^{M} \alpha_i = 1 \tag{3.6}$$

ensures that the total probability that x_j came from any of the M distributions sums to one.

Put forth in the same terms as were used in section 3.1:

$$Y = F\left(X\right)$$

where the goal is to recover the value of the unknown class membership Y, from the data X using mapping F. Y contains information as to which of the M Gaussian distributions datum x_j came from. The problem is complicated by the fact that the mapping function, F, is also unknown, since the parameters θ which define each of the M Gaussian distributions have yet to be determined. Finally, as with K-Means, there still lies to matter of choosing a number of clusters (Gaussian distributions), M, to represent the data.

Therefore given known data X, and the assumption that the data are derived from a Gaussian Mixture Model, the following variables will need to be solved for:

- \boldsymbol{Y} the class memberships. This contains information as to which of the M Gaussian distributions each datum is most likely to come from.
- $\boldsymbol{\theta}$ the defining parameters for each of the *M* Gaussian distributions. $\theta_i = [\mu_i, \Sigma_i]$ i = 1..M.
- α the weight given to each distribution, 1..*M*. This value dictates the importance of each Gaussian distribution in representing the data.

In the following sections the EM method will be derived to solve for these quantities.

3.3.2 Maximum-Likelihood

Since the data are independent and identically distributed (i.i.d.), the likelihood of observing the entire data X given parameters θ is the product of the likelihood of observing each datum independently:

$$l(X|\theta) = \prod_{j=1}^{N} l(x_j|\theta) = \prod_{j=1}^{N} \sum_{i=1}^{M} \alpha_i g_i(x_j|\theta_i)$$
(3.7)

The goal of a Maximum Likelihood algorithm is to obtain an estimate for the parameters of each Gaussian distribution, $\theta = [\theta_1, \theta_2, ..., \theta_M]$ and $\alpha = [\alpha_1, \alpha_2, ..., \alpha_M]$, where $\theta_i = (\mu_i, \Sigma_i)$ such that they maximize the likelihood function, $l(X|\theta)$. In other words, to discover the most likely parameters $[\theta, \alpha]$ for the generative distribution $l(X|\theta)$ which explains the data X. Because $l(X|\theta)$ is a product (Eq. 3.7), the maximization can be simplified by taking its logarithm³ and maximizing this instead:

³Note: For all discussion of the EM Method, *log* refers to the natural logarithm.

$$log\left[l\left(X|\theta\right)\right] = \sum_{j=1}^{N} log\left[\sum_{i=1}^{M} \alpha_{i} g_{i}\left(x_{j}|\theta_{i}\right)\right]$$
(3.8)

which will be represented as $L(\theta)$ since the log-likelihood is only a function of the parameters θ .

3.3.3 E-Step

The E-Step, or Expectation step, applies the expectation value along with some other mathematical tricks to rewrite the log-likelihood, $L(\theta)$, in terms that can be more easily maximized.

Applying Bayes theorem (Hastie et al. [2001]), it can be shown that since the likelihood that the data came from any of the M Gaussian distributions is equal to unity

$$\sum_{i=1}^{M} l(y_i | X, \theta) = \mathbf{1}$$
(3.9)

the log-likelihood can also written as

$$L(\theta) = \log\left(l(X|\theta)\right) = \log\left(\sum_{i=1}^{M} l(X, Y|\theta)\right)$$
(3.10)

Multiplying this by unity (Eq. 3.9) gives

$$L(\theta) = \log\left(\sum_{i=1}^{M} l(y_i|X, \theta^p) \frac{l(X, Y|\theta)}{l(y_i|X, \theta^p)}\right)$$
(3.11)

where θ^p is defined as the current best estimate of the parameters at the p^{th} iteration. Jensen's inequality (Hastie et al. [2001]) can be invoked to rewrite this as:

$$L(\theta) \ge \sum_{i=1}^{M} l(y_i | X, \theta^p) \log\left(\frac{l(X, Y | \theta)}{l(y_i | X, \theta^p)}\right)$$
(3.12)

Finally, splitting up the logarithmic function gives two terms:

$$L(\theta) \geq \sum_{i=1}^{M} l(y_i|X,\theta^p) \log (l(X,Y|\theta)) - \sum_{i=1}^{M} l(y_i|X,\theta^p) \log (l(y_i|X,\theta^p))$$

$$\geq Q(\theta|\theta^p) - R(\theta^p|\theta^p)$$
(3.13)

Since the goal is to maximize this expression (by taking the derivative) with respect to θ , only $Q(\theta|\theta^p)$ is of interest since $R(\theta^p|\theta^p)$ is not a function of θ . Therefore only the first term in Eq. 3.13 is of interest, and can be explicitly written as:

$$Q_j(\theta|\theta^p) = \sum_{i=1}^M \gamma_{ij}^p log\left(l\left(x_j, y_j|\theta\right)\right)$$
(3.14)

where

$$\gamma_{ij}^p = l\left(y_j = i | X = x_j, \theta^p\right) \tag{3.15}$$

Explicitly this is evaluating the following:

$$\gamma_{ij}^{p} = \frac{\frac{\alpha_{i}^{p}}{(2\pi)^{\frac{K}{2}} \det\left(\Sigma_{i}^{p}\right)^{\frac{1}{2}}} exp\left[\frac{-1}{2} \left(x_{j} - \mu_{i}^{p}\right)^{T} \left(\Sigma_{i}^{p}\right)^{-1} \left(x_{j} - \mu_{i}^{p}\right)\right]}{\sum_{k=1}^{M} \frac{\alpha_{k}^{p}}{(2\pi)^{\frac{K}{2}} \det\left(\Sigma_{k}^{p}\right)^{\frac{1}{2}}} exp\left[\frac{-1}{2} \left(x_{j} - \mu_{k}^{p}\right)^{T} \left(\Sigma_{k}^{p}\right)^{-1} \left(x_{j} - \mu_{k}^{p}\right)\right]} \quad (3.16)$$

In plain text, γ_{ij}^p can be thought of as the best guess as to the probability that x_i belongs to component *i* at iteration *p*.

Plugging in Eq. 3.5 and simplifying, the expected value of the log-likelihood of the datum x_j coming from distribution *i* given parameters θ^p becomes maximized by:

$$Q_{j}(\theta|\theta^{p}) = \sum_{i=1}^{M} \gamma_{ij}^{p} \left(\log \alpha_{i} - \frac{1}{2} \log |\Sigma_{i}| - \frac{1}{2} (x_{j} - \mu_{i})^{T} (\Sigma_{i})^{-1} (x_{j} - \mu_{i}) \right)$$
(3.17)

Summing over all data, the expected value of the complete data log-likelihood is maximized by:

$$Q(\theta|\theta^{p}) = \sum_{j=1}^{N} \sum_{i=1}^{M} \gamma_{ij}^{p} \left(\log \alpha_{i} - \frac{1}{2} \log |\Sigma_{i}| - \frac{1}{2} (x_{j} - \mu_{i})^{T} (\Sigma_{i})^{-1} (x_{j} - \mu_{i}) \right)$$
(3.18)

This can be simplified further by defining the following as the number of data in each cluster at the end of the p^{th} iteration:

$$n_i^p = \sum_{j=1}^N \gamma_{ij}^p \tag{3.19}$$

So that Eq. 3.18 now becomes:

$$Q(\theta|\theta^{p}) = \sum_{i=1}^{M} n_{i}^{p} \left(\log \alpha_{i} - \frac{1}{2} \log |\Sigma_{i}| \right) - \frac{1}{2} \sum_{j=1}^{N} \sum_{i=1}^{M} \gamma_{ij}^{p} \left(x_{j} - \mu_{i} \right)^{T} \left(\Sigma_{i} \right)^{-1} \left(x_{j} - \mu_{i} \right)$$
(3.20)

3.3.4 M-Step

Now that the log-likelihood has been rewritten in more manageable terms, the function in Eq. 3.20 will be maximized with respect the unknown variables: μ_i, Σ_i and α_i , for i = 1..M. This will be done by taking derivatives and setting them equal to zero.

Mixing Weights α

The first set of variables which will be solved for are the mixing weights, α . Formally, the following will be solved:

$$\underset{\theta}{\operatorname{argmax}} Q\left(\theta|\theta^{p}\right) \quad s.t. \quad \sum_{i=1}^{M} \alpha_{i} = 1 \tag{3.21}$$

Using a Lagrangian to solve for the weights gives

$$J(\alpha, \lambda) = Q(\theta|\theta^p) + \lambda \left(\sum_{i=1}^{M} \alpha_i - 1\right)$$
(3.22)

Taking the derivative of this expression with respect to the weights

$$\frac{\partial J\left(\alpha,\lambda\right)}{\partial\alpha_{i}} = \frac{n_{i}^{p}}{\alpha_{i}} + \lambda = 0 \tag{3.23}$$

Summing over all the components, i = 1..M, gives

$$\lambda = -\sum_{i=1}^{M} n_i^p \tag{3.24}$$

Therefore

$$\alpha_i^{p+1} = \frac{n_i^p}{\sum_{i=1}^M n_i^p} = \frac{n_i^p}{N}$$
(3.25)

which intuitively makes sense: the weight (importance) prescribed to each Gaussian distribution is equal to the number of data which are represented by each distribution normalized by the total number of data.

Mean Values μ

Following the same procedure, the mean values can be solved for by taking the derivative of Eq. 3.20, this time with respect to μ_i and setting it equal to zero:

$$\frac{\partial Q\left(\theta|\theta^{p}\right)}{\partial \mu_{i}} = \left(\Sigma_{i}^{p}\right)^{-1} \left(\sum_{j=1}^{N} \gamma_{ij}^{p} x_{j} - n_{i}^{p}\right) = 0 \qquad (3.26)$$

Rearranging gives

$$\mu_i^{p+1} = \frac{1}{n_i^p} \sum_{j=1}^N \gamma_{ij}^p x_j \tag{3.27}$$

42

which, for the i^{th} Gaussian distribution, in words is: "the sum of all data (x_j) weighted by their probability of belonging to the i^{th} distribution (γ_{ij}) , then normalized by the number of data represented by the i^{th} distribution (n_i) ".

Covariances Σ

Finally, once more to solve for the covariances, the derivative of Eq.3.20 is taken with respect to Σ_i , and set equal to zero:

$$\frac{\partial Q\left(\theta|\theta^{p}\right)}{\partial \Sigma_{i}} = \frac{-1}{2}n_{i}^{p}\Sigma_{i}^{-1} + \frac{1}{2}\sum_{j=1}^{N}\gamma_{ij}^{p}\Sigma_{i}^{-1}\left(x_{j}-\mu_{i}^{p}\right)\left(x_{j}-\mu_{i}^{p}\right)^{T}\Sigma_{i}^{-1} = 0 \quad (3.28)$$

Rearranging gives

$$\Sigma_{i}^{p+1} = \frac{1}{n_{i}^{p}} \sum_{j=1}^{N} \gamma_{ij}^{p} \left(x_{j} - \mu_{i}^{p} \right) \left(x_{j} - \mu_{i}^{p} \right)^{T}$$
(3.29)

which again, for the i^{th} Gaussian distribution, in words is: "the sum of the squared difference between each datum and the i^{th} mean value (μ_i) , weighted by their probability of belonging to the i^{th} distribution (γ_{ij}) , and then normalized by the number of data represented by the i^{th} distribution (n_i) ".

3.3.5 Convergence

These two steps, the E-Step and the M-Step, are iterated repeatedly until convergence is reached. The typical stopping criterion for this being

$$\left|L\left(\theta^{p+1}\right) - L\left(\theta^{p}\right)\right| < \delta \tag{3.30}$$

In other words, when the difference between the last iteration's log-likelihood and the current iteration's log-likelihood is less than some predefined threshold value, δ , stop the algorithm. Since the M-Step acts at each iteration to

maximize $Q(\theta|\theta^p)$, a new set of parameters θ^{p+1} can always be chosen such that $Q(\theta^{p+1}|\theta^p) \ge Q(\theta|\theta^p)$. Given Eq. 3.13, it can be shown that since $L(X|\theta^{p+1}) \ge Q(\theta^{p+1}|\theta^p) - R(\theta^p|\theta^p) \ge Q(\theta^p|\theta^p) - R(\theta^p|\theta^p) = L(X|\theta^p)$ (3.31)

this implies that $L(\theta^{p+1}) \geq L(\theta^p)$, and thus the log-likelihood never decreases.

3.3.6 Numerical Example: EM Algorithm

The EM algorithm will now be applied to our toy problem:

j =	1	2	3	4	5	6	7
	-0	-0	-0-	-0		_0_	-0
$x_j^{=}$	1	2	3	4	6	7	8

Figure 3.3: Simple 1D data example

As with K-Means, two clusters will be assumed, however for the EM algorithm this implies that there exists a mixture of two one-dimensional Gaussian distributions. To begin, the problem will be initialized with the following parameters:

$$\begin{array}{rclrcl}
\mu_1^0 &=& 0.0 & \mu_2^0 &=& 9.0 \\
\Sigma_1^0 &=& 1.0 & \Sigma_2^0 &=& 1.0 \\
\alpha_1^0 &=& 0.5 & \alpha_2^0 &=& 0.5
\end{array}$$
(3.32)

Note that this follows the same notation as the previous sub-section, therefore θ_i^p are the parameters for the i^{th} distribution at the p^{th} iteration.

3.3.6.1 E-Step

As per Eq. 3.16 γ_{1j}^1 will be calculated as:

$$\gamma_{1j}^{1} = \frac{\frac{\alpha_{1}^{0}exp\left[\frac{-1}{2}\left(x_{j}-\mu_{1}^{0}\right)^{T}\left(\Sigma_{1}^{0}\right)^{-1}\left(x_{j}-\mu_{1}^{0}\right)\right]}{\left(2\pi\right)^{\frac{K}{2}}\left|\Sigma_{1}^{0}\right|^{\frac{1}{2}}}}{\frac{\alpha_{1}^{0}exp\left[\frac{-1}{2}\left(x_{j}-\mu_{1}^{0}\right)^{T}\left(\Sigma_{1}^{0}\right)^{-1}\left(x_{j}-\mu_{1}^{0}\right)\right]}{\left(2\pi\right)^{\frac{K}{2}}\left|\Sigma_{1}^{0}\right|^{\frac{1}{2}}} + \frac{\alpha_{2}^{0}exp\left[\frac{-1}{2}\left(x_{j}-\mu_{2}^{0}\right)^{T}\left(\Sigma_{2}^{0}\right)^{-1}\left(x_{j}-\mu_{2}^{0}\right)\right]}{\left(2\pi\right)^{\frac{K}{2}}\left|\Sigma_{1}^{0}\right|^{\frac{1}{2}}}$$

$$(3.33)$$

Plugging in values:

$$\gamma_{1j}^{1} = \frac{\frac{(0.5)exp\left[\frac{-1}{2}(x_{j}-(0.0))^{T}(1.0)^{-1}(x_{j}-(0.0))\right]}{(2\pi)^{\frac{1}{2}}|1.0|^{\frac{1}{2}}}{\frac{(0.5)exp\left[\frac{-1}{2}(x_{j}-(0.0))^{T}(1.0)^{-1}(x_{j}-(0.0))\right]}{(2\pi)^{\frac{1}{2}}|1.0|^{\frac{1}{2}}} + \frac{(0.5)exp\left[\frac{-1}{2}(x_{j}-(9.0))^{T}(1.0)^{-1}(x_{j}-(9.0))\right]}{(2\pi)^{\frac{1}{2}}|1.0|^{\frac{1}{2}}}$$

$$(3.34)$$

Evaluating this for all x_j , j = 1...N, the following results are obtained:

j	x_j	γ_{1j}^1	γ_{2j}^1
1	1	1.0	2.088×10^{-14}
2	2	1.0	1.6919×10^{-10}
3	3	1.0	1.371×10^{-6}
4	4	0.98901	0.010987
5	6	1.371×10^{-6}	1.0
6	7	1.6919×10^{-10}	1.0
7	8	2.088×10^{-14}	1.0

Table 3.2: γ^1_{ij} Values

Recall that γ_{ij} denotes the probability of datum x_j belonging to the i^{th} Gaussian distribution, and that (to mathematical precision) $\sum_{i=1}^{M} \gamma_{ij} = 1$. Due to the simplicity of this problem, one can see that by the first iteration the EM algorithm has already assigned data 1-4 to the first distribution, and data 5-7 to the second.

3.3.6.2 M-Step

Now that γ_{ij}^1 has been calculated, a new set of parameters, μ_1^1 , Σ_1^1 , α_1^1 and μ_2^1 , Σ_2^1 , α_2^1 can be estimated as per Eqs. 3.25, 3.27 and 3.29:

$$\mu_1^1 = \frac{\sum_{j=1}^N \gamma_{1j}^1 x_j}{\sum_{j=1}^N \gamma_{1j}^1} \tag{3.35}$$

$$\Sigma_{1}^{1} = \frac{\sum_{j=1}^{N} \gamma_{1j}^{1} \left(x_{j} - \mu_{1}^{1}\right) \left(x_{j} - \mu_{1}^{1}\right)^{T}}{\sum_{j=1}^{N} \gamma_{1j}^{1}}$$
(3.36)

45

$$\alpha_1^1 = \frac{\sum_{j=1}^N \gamma_{1j}^1}{N} \tag{3.37}$$

Evaluating these expressions, as well as those for the second Gaussian distribution, for all iterations, we obtain these values:

Iteration (p)	μ_1^p	Σ_1^p	α_1^p	μ_2^p	Σ_2^p	α_2^p
$\operatorname{K-Means}$	2.50	1.25	-	7.00	1.00	-
0	0	1	0.5	9	1	0.5
1	2.4959	7.4766	0.5699	6.9891	4.7409	0.4301
2	3.0864	4.0486	0.5921	6.3769	4.7409	0.4079
3	2.9743	3.2647	0.5918	6.5367	3.2027	0.4082
4	2.7912	2.5497	0.5875	6.7608	2.3991	0.4125
5	2.6427	1.8692	0.5829	6.9242	1.6556	0.4171
6	2.5606	1.4777	0.5757	6.9970	1.0472	0.4211
7	2.5289	1.3471	0.5743	7.0064	0.7413	0.4243
8	2.5194	1.3140	0.5739	7.0047	0.6768	0.4257
9	2.5168	1.3060	0.5738	7.0038	0.6724	0.4261
10	2.5162	1.3040	0.5738	7.0035	0.6726	0.4262
11	2.5161	1.3035	0.5738	7.0034	0.6728	0.4262
12	2.5160	1.3034	0.5738	7.0034	0.6729	0.4262
13	2.5160	1.3034	0.5738	7.0034	0.6729	0.4262
14	2.5160	1.3034	0.5738	7.0034	0.6729	0.4262

 Table 3.3: Parameter Values

where the results from the K-Means algorithm have been marked in blue, and the values to which EM parameters converged are marked in red. The mean values are almost identical, while the variances are slightly different. This is due to the fundamental difference between the two algorithms: K-Means fits each data to a single cluster, whereas the EM algorithm will fit every datum to every cluster, and thus the variance will be slightly larger to accommodate points which are further away. Comparing figure 3.2 to 3.4 below, it would appear that the EM method has better fit the distribution of data. Additionally, as a metric of the weighting values, α , consider that $\frac{4}{7}$ data gives a weight of 0.5714, and $\frac{3}{7}$ data a weight of 0.4286.



Figure 3.4: Toy problem: EM algorithm results

3.3.6.3 Convergence

Given the values in Table 3.3 the log-likelihood can be calculated at each iteration. As was previously asserted, the values of $Q(\theta|\theta^p)$ are always greater or equal to the previous values (see Table 3.4).

Iteration (p)	$Q\left(heta heta^{p} ight)$	$L\left(heta ight)$
0	-26.9015	-2.1121
1	-13.2172	-2.5044
2	-12.0333	-2.7148
3	-11.0431	-2.6589
4	-9.8202	-2.5115
5	-8.7451	-2.3682
6	-8.3068	-2.3128
7	-8.2328	-2.3104
8	-8.2201	-2.3113
9	-8.2174	-2.3116
10	-8.2168	-2.3116
11	-8.2167	-2.3117
12	-8.2166	-2.3117
13	-8.2166	-2.3117
14	-8.2166	-2.3117

Table 3.4: Convergence of the EM Method

3.3.7 Practical Issues

3.3.7.1 Number of Clusters

As was stated in section 3.3.1, one of the drawbacks of classification techniques such as the *Expectation-Maximization* method is that they require

the user to specify the desired number of clusters, M. In many situations there is no way of knowing the exact number of distributions contributing to the resulting data. Since relocation methods such as *K*-Means and the EM method can be very unstable, the choice of M can be very important to the result. Since Bayesian statistics are used to derive the result, a common metric for the goodness-of-fit is the AIC, or Akaike Information Criterion, defined as the following:

$$AIC = -2L\left(\theta\right) + 2k \tag{3.38}$$

where $L(\theta)$ is the log-likelihood of the resulting data, X, and associated memberships, Y, given the estimated parameters θ , k is the number of parameters being estimated, and N is the number of data. Using the AIC, the question of how many clusters becomes the following optimization problem:

$$\mathop{argmin}_{M} AIC\left(M,\theta\right)$$

Plotting the AIC gives a curve similar to the following



Figure 3.5: Plot of number of clusters vs. Akaike information criterion

where the first local minimum⁴ (indicated by a \diamond) suggests strong evidence for the corresponding model. A further refinement on the *AIC* for cases when the number of parameters is large relative to the number of data is the *AICc*:

$$AICc = AIC + \frac{2k(k+1)}{N-k-1} = -2L(\theta) + 2k + \frac{2k(k+1)}{N-k-1}$$
(3.39)

Other criteria such as the BIC, or Bayes Information Criterion

$$BIC = -2L(\theta) + klog(N) \tag{3.40}$$

exist, and although debate exists as to which should be preferred as a metric of the optimal model, the literature seems to slightly favor the AIC, and thus it has been used in this work.

3.3.7.2 Initialization

Because the *Expectation-Maximization* method is unstable, it is sensitive to the initial conditions. In order to mitigate these effects, the algorithm is initiated by choosing M data at random to assign as the initial mean values, with diagonal covariances such that Σ_{ij} for i = j is var(X(:,j)), and for $i \neq j$ is 0.

The algorithm is started this way a number of times, with a new random set of initial mean values chosen each time, and the repeat with the highest log-likelihood is kept as the best run.

As with any iterative algorithm it is also necessary to define a maximum number of iterations so that it does not get stuck in an infinite loop. Should the algorithm fail to converge due to ill-conditioning, regularization is sometimes necessary.

3.3.7.3 Regularization

Since the M-Step of the EM method requires evaluating Σ^{-1} it is important that Σ is in fact invertible. Many factors can potentially lead to an ill-conditioned covariance matrix, such as:

 $^{^4 \}rm Some$ authors have formulated the BIC as the negative of Eq.3.38, and thus the first local maximum is used instead.

- If the number of data N is small relative to the dimensionality M of the data
- If the data are highly correlated
- Too many clusters are being used

Should one or more of these difficulties be present, it is possible that the algorithm will converge to a local solution which has an ill-conditioned covariance matrix. In order to avoid this, one strategy is to regularize the covariance matrix by adding a small ($\sim 10^{-6}$) positive value to the diagonal components. Another less-optimal strategy is to use shared covariance matrices; in other words all the covariance matrices are the pooled estimate of the entire data-set. This restricts the ability of the algorithm to model complex distributions, but it can resolve the ill-conditioning problem.

3.4 Geophysical Application

Despite only presenting the Expectation-Maximization algorithm in any reasonable detail, other classification schemes were also attempted. In particular Self Organizing Maps, or SOMs (Kohonen [1990], Vesanto et al. [2000]), were investigated as a potential classification algorithm, however in all initial testing⁵ the EM algorithm was able to correctly classify the model at least as well as SOMs, if not better. This is likely due in part to the fact that the EM algorithm assumes normally distributed point clouds of data values, which is also the assumption made for the distribution of physical property values in the earth. It is possible that with further investigation, SOMs might present an alternative effective means of classifying downhole data, without the need for such assumptions (Fraser and Hodgkinson [2009]).

Though some of the finer elements of the statistical classification discussed in this chapter may be difficult to grasp, the important point is that classification has been presented as a viable means to process the vast amount of information supplied by downhole physical property logs. Not only is classification able to leverage the high vertical resolution of the physical property

⁵Prior to applying any classification algorithms to downhole data, various algorithms (K-Means, Fuzzy Logic C-Means, SOMs, and EM) were tested on a number of synthetic distributions of data with varying degrees of difficulty (more overlap in clusters, wider distribution spreads etc).

logs, but also to integrate the information from various physical properties into one coherent log of "rock type" with depth. Associated with this "rock type" log are a set of parameters which define each geophysical unit. For each physical property measured in the classified logs, a mean value and associated standard deviation are supplied for each unique rock type the algorithm finds.

The statistical classification of various defining parameters of the earth (be they geological, geochemical or geophysical) has long been applied to help understand distributions of natural resources (Journel and Huijbregts [1978], Isaaks and Srivastava [1989], Matheron [1963]). Over the last couple decades, geostatistics has been applied to assist in understanding the relation and distribution of multiple physical properties from downhole logs (MacMahon et al. [2002], S.E. MacMahon [2002]). Traditionally the results of such analysis have been interpreted on their own, applying methods such as kriging to develop models of the subsurface. Only recently has research begun to be dedicated to the incorporation of geostatistics to constraining geophysical inversion (Wang and Yang [2011], Hermans et al. [2011]). Currently, the majority of efforts are taking place in petroleum, environmental, and engineering geophysics, however applications for mineral exploration are becoming more common.

In the following chapter a basic introduction of the theory of geophysical inversion will be presented so that the means of constraining inversion can be better understood. With this knowledge, an explanation of a methodology for incorporating classification results into the creation of constraints will be better understood.

Chapter 4

Geophysical Inversion

4.1 Introduction

In geophysics, measurements of the earth are taken so that a set of desired parameters might be inferred about the underlying volume of earth, such as location or size of a body, or some distribution of physical properties. The difficulty lies in the fact that it is not typically possible to directly measure these parameters, and therefore an experiment must be set up in which some knowledge of the physics involved is assumed.

In general, such an experiment will involve an input of energy to the earth, followed by measurements of the energy output. Specifics such as what kind of energy is input, how the energy propagates through the system, or how best to measure the energy output are all dependent on the system being studied and the parameters one hopes to estimate. Assuming these details are correctly chosen, the goal is to arrive at a representative model which is a best estimate of the parameters in question.

4.1.1 Forward Problem

In order to determine whether or not the optimal model has been chosen, a basic but important metric is whether the selected model is able to reproduce the observed data. This requires the ability to simulate measurements of an arbitrary model - termed the *forward problem*. In mathematical terms, this is equivalent to applying a functional F to a vector m in model space to arrive at a vector d in data space (see figure 4.1). Practically, this is the act of taking measurements, since it is our measurements which give us information about the model in question (wherein the physics involved can be thought of as the functional).





Figure 4.1: Diagram illustrating forward modeling

4.1.2 Inverse Problem

As one might expect, the *inverse problem* is just the opposite of the *for-ward problem*. Attempting to recover parameters from experimental measurements is in effect attempting to solve an *inverse problem* - hence *inversion*. Here, we have a vector in data space d, and we would like to map it back to a vector in model space m (see figure 4.2).



Figure 4.2: Diagram illustrating inversion

This problem is infinitely more difficult than the *forward problem*, and before we attempt to solve it we must ask ourselves some important questions:

- Does a solution exist?
- If a solution does exist, is it unique?
- If the solution exists but is not unique, are there properties that are uniquely determined?

Without going into too many details, it can be stated that the *inverse problem* is fundamentally non-unique (Oldenburg and Li [2005]) since there exist an infinite number of models which could reproduce the data. Therefore solving these problems requires a little extra thought. In the next section a brief outline of the mathematical formulation for a typical geophysical inversion will be presented, including the means by which the problem of non-uniqueness is tackled.

4.2 Mathematical Formulation

As was mentioned in the previous section, the main difficulty with *inverse* problems is that they are fundamentally non-unique. In order to solve them, we must therefore constrain the number of possible solutions (Oldenburg and Pratt [2007]). There are a number of different ways this can be done, however for our applications we will be applying *Tikhonov Regularization* of the form:

$$\phi = \phi_d + \beta \phi_m \tag{4.1}$$

where ϕ_m is the objective function, ϕ_d is the data misfit, and β is the Tikhonov parameter. Below we will work to build up the mathematical formulation of each of these components, and to give practical justification for each of the constraints.

4.2.1 Data Misfit

The most basic criteria for a correct model is that it should be able to reproduce the observed data to within a reasonable threshold. Consider the following geophysical data:

$$d_j = (g_j, m) \qquad j = 1..N$$
 (4.2)

where d_j is the projection of the j^{th} basis vector, g_j , onto the model, m. Or, equivalently in matrix form:

$$d^{obs} = Gm \tag{4.3}$$

Then if the predicted model is m^{pred} , the predicted data will be

$$d^{pred} = Gm^{pred} \tag{4.4}$$

and we would like to minimize the difference between the two:

$$\phi_d = \| d^{pred} - d^{obs} \|_2 = \| Gm^{pred} - d^{obs} \|_2$$
(4.5)

In reality the goal is not really to minimize the data misfit. If one considers noisy data:

$$d^{obs} = Gm + \epsilon \tag{4.6}$$

where ϵ is Gaussian noise, then one can re-represent the data misfit as:

$$\phi_d = \parallel W_d(Gm^{pred} - d^{obs}) \parallel_2$$

where

$$W_d = diag\left(\frac{1}{\sigma_i}...\frac{1}{\sigma_N}\right)$$

where σ is the standard deviation, N is the number of data, and W_d is referred to as a *weighting function*. The effect of the weighting function is to normalize the noise on each datum, therefore the desired difference between the predicted data and the noisy data is simply the number of data

$$\phi_d \equiv N \tag{4.8}$$

4.2.2 Objective Function

The objective function specifies the desired behavior of the model. In a typical geophysical inversion there will be two main components to the objective function: smallness and smoothness.

Smallness, or smallest deviatoric model, ensures that the recovered model is similar to a desired reference model, m_{ref} . Similarly to the data misfit, the goal here is to minimize the difference between the recovered model and the reference model:

$$\phi_s = \| W_s (m - m_{ref}) \|_2 \tag{4.9}$$

(4.7)
where W_s is again a weighting function which allows the user to control which elements of the recovered model should most resemble the reference model.

The other component, smoothness, can be applied in any spatial direction. This component helps to prevent large discontinuities from appearing in the model by enforcing that the recovered model be smooth in a given direction. This is accomplished by minimizing the derivative of the difference between the recovered model and the reference model:

$$\phi_{x_i} = \| W_{x_i} \frac{d}{dx_i} \left(m - m_{ref} \right) \|_2 \tag{4.10}$$

where $\frac{d}{dx_i}$ is the derivative in the x_i direction, and as always W_{x_i} is a weighting function. Therefore in three dimensions one could easily imagine having three smoothness components in the objective function. Adding it all together along with weighting parameters α_i , we arrive at the following expression for the model objective function:

$$\phi_{m} = \alpha_{s} \| W_{s} (m - m_{ref}) \|_{2}
+ \alpha_{x} \| W_{x} \frac{d}{dx} (m - m_{ref}) \|_{2}
+ \alpha_{y} \| W_{y} \frac{d}{dy} (m - m_{ref}) \|_{2}
+ \alpha_{z} \| W_{z} \frac{d}{dz} (m - m_{ref}) \|_{2}$$
(4.11)

The weighting parameters, α_i , are constants which allow the user to specify the importance of a given constraint (*ie*: smallness or smoothness). For example, if a layered earth model is suspected, the user might put a larger emphasis on the smoothness of the horizontal components (x and y) than on the vertical (z), since discontinuities would be expected between vertical layers. This would be accomplished by making $\alpha_z < \alpha_x, \alpha_y$.

To simplify this, the weighting functions, weighting parameters and derivatives can all be collapsed into one large matrix, W_m , resulting in the more common notation:

$$\phi_m = \| W_m (m - m_{ref}) \|_2 \tag{4.12}$$

4.2.3 Solving

Combining the results from Sections 4.2.1 and 4.2.2, and rewriting Eq. 4.1, we arrive at the following expression to minimize:

$$\phi = \| W_d(Gm^{pred} - d^{obs}) \|_2 + \beta \| W_m(m - m_{ref}) \|_2$$
(4.13)

The minimization is carried out by taking the gradient with respect to m and setting it equal to zero. After a little linear algebra we arrive at the following:

$$\left(G^T W_d^T W_d G + \beta W_m^T W_m\right) m = G^T W_d^T W_d d^{obs} + \beta W_m^T W_m m_{ref} \qquad (4.14)$$

The expression is then solved for a number of different β values in a line search to find β^* such that $\phi_d = \phi_d^* = N$. This produces what is referred to as a Tikhonov curve, or L-curve, seen below.



Figure 4.3: Diagram of the Tikhonov curve used for solving geophysical inverse problem

It should be noted that for large problems this methodology is not always preferable, or possible, due to computational limitations and difficulties. A suite of techniques exist to tackle these problems, and for more information the reader is directed to Oldenburg and Pratt [2007], Oldenburg and Li [2005].

4.2.4 Control Parameters

In order to handle a wide array of different problems, Eq. 4.14 is equipped with a number of tunable parameters which allow the user to apply prior information. In this section they will be discussed with reference to both their impact on an inversion, as well as their physical interpretation. In 3-dimensions, the 12 parameters fall into three main groups:

Weighting Functions	Weighting Parameters	Reference Models
$W_x \\ W_y$	$lpha_x \ lpha_y$	${m_{ref} \over m^{min}}$
W_z W_s	α_z α_s	m^{max}
W_d		

Table 4.1: Tunable Parameters in UBC GIF 3D Inversion Code

4.2.4.1 Weighting Functions

Weighting functions, W_i , allow the user to distribute confidence estimates unevenly across the spatial extent of the model space. Since each weighting function is a matrix of size equal to the number of cells in model space, a different weighting can be applied to each cell. In this way it is possible to incorporate more inhomogeneous structures into the recovered model. Consider the following example:

Suppose we are trying to recover a line from point A to point B, and we have a priori knowledge of the value of the model at two points, x_1 and x_2 (see figure 4.4 below).



Figure 4.4: Unknown 1D model which connects points A and B, with prior knowledge of the values at points x_1 and x_2

One way to incorporate this information into the recovered model is to apply a weighting function of the form:



Figure 4.5: Weighting function, W_x used to emphasize known points in the model. Larger values will be more penalized for not resembling the reference model.

In this way, at points x_1 and x_2 the model will be heavily biased towards the reference model, m_{ref} , which might look something like:



Figure 4.6: Example of 1D reference model, m_{ref}

while the rest of the model will be less heavily impacted. This can be generalized to apply to any of the weighting functions listed in Table 4.1, resulting in a wide range of spatial variability in the recovered model.

The one weighting function which might be considered different than the rest (though it serves the same purpose) is the data misfit weighting function, W_d . As was suggested in Eq. 4.7, the data misfit weighting function is typically taken to be a diagonal matrix of one over the standard deviation:

$$W_d = diag\left(\frac{1}{\sigma_i}...\frac{1}{\sigma_N}\right) \tag{4.15}$$

Applying this weighting to the data misfit has the effect of normalizing the data misfit by the standard deviation of each datum. When this is done, the desired misfit is

$$\phi_d^* = N \tag{4.16}$$

4.2.4.2 Reference Models

The reference model, m_{ref} , plays a crucial role in the inversion process. It is in the reference model that the user is able to specify an estimate of the expected recovered model. Referring back to Section 4.2, it is evident that the reference model appears in many of the constraint equations, and therefore the choice of reference model can have large implications for the recovered model.

4.2.4.3 Bounds

The other two constraints, m^{min} and m^{max} , are optional additional models that can be included to impose upper and lower bounds on the value of each cell in the resulting discretized model. Mathematically, this additional constraint modifies Eq. 4.1 so that it is now

$$\begin{array}{ll} minimize & \phi = \phi_d + \beta \phi_m \\ s.t. & m^{min} \leqslant m \leqslant m^{max} \end{array}$$
(4.17)

This kind of a constraint can be very useful when some information is available about the expected geology being represented in the model, or for certain physical properties which have known physically justified constraints (ie: positivity).

4.3 DCIP2D

This thesis employs DCIP2D to carry out the modeling of DC resistivity data. Developed at the University of British Columbia Geophysical Inversion Facility, the software package is able to both simulate DC resistivity data through forward modeling, as well as recover two-dimensional resistivity models through the inversion of surface data.

In order to provide adequate flexibility to handle a variety of different geophysical targets, DCIP2D is equipped with a number optional input parameters. Of these parameters, only a handful were applied in this thesis. These include the mesh file and the data file (both of which are self explanatory), as well as an initial/reference model, alpha values, weight files, and bounds files.

The initial model and reference model, as mentioned in section 4.2.4, allow the user to specify an estimate of the expected model. If the recovered model, m, deviates from this reference model, m_{ref} , it is penalized according to the model objection function (equation 4.11). Additionally, whether or not m_{ref} appears in the derivative terms (equation 4.10) of the objective function is specified by USE_MREF. When set as FALSE, m_{ref} will not appear in these terms, and thus the inversion is not penalized if the derivative of the difference between the reference model and the recovered model is large.

This can further be refined by specifying the alpha values used in the objective function (equation 4.11). These values $(\alpha_s, \alpha_x \text{ and } \alpha_z)$ help to control the importance of each term in the objective function. By making α_s small, the emphasis is placed on recovering a smooth model rather than a model which is similar to the reference model, and vice versa.

Weighting files can be applied to enforce smooth or sharp boundaries in specific locations in the model. One typical application of weighting files with DCIP inversions is to mitigate noise close to the surface in the recovered model. This noise is characteristic of a DCIP inversion, and often occurs close to the electrode locations in the model. By using weight files to horizontally smooth the surface, these effects can be minimized without adversely affecting the entire model.

Finally, bounds files can be applied to constrain the recovered resistivity values to within some predefined boundaries. This is particularly useful when accurate information is available to limit the range of possible values within the model. Furthermore, since the bounds are applied on a cell by cell basis, known general structure can be incorporated into the inversion without the need to lock in specific values or lithological boundaries. In the following chapter the topics discussed in the previous three chapters will be combined to develop an iterative inversion technique which incorporates the information provided by downhole physical property logging into the inversion of surface geophysics. If the reader is interested in more information on the application of the DCIP2D software package, they are directed to the user manual (Uni [2011]).

Chapter 5

Iterative Inversion Technique

Given the tools provided in the last three chapters, an attempt was made to maximize the incorporation of valuable information acquired from downhole physical property logs into the inversion of surface geophysical data. The resulting multi-step procedure is explained below, followed by a set of simple illustrative examples.

5.1 Methodology

The methodology explained below (figure 5.1) aims to apply the information from downhole logs to constrain geophysical inversion, and to do so with minimal bias from the user. This is achieved via the statistical classification of downhole physical property logs discussed in chapter 3, and a departure from soft constraints such as reference models, to the harder and more flexible bounds constraints. By iteratively updating the bounds based on the results from classification, a final model can be achieved which is both accurate and reliable, with minimal information provided by the user, and maximum application of the physical property logs.

To demonstrate this procedure, synthetic modeling was performed. The goal of synthetic modeling was to simulate real data using a derived geophysical model, and apply the suggested methodology in an attempt to recover the model from the data. This demonstration allows for better understanding of the entire process since the models and data are well known and controlled, facilitating the evaluation of the recovered models.



Figure 5.1: Methodology flow chart

The synthetic modeling began with the creation of simple geological models. Each geological model was then translated into a set of physical property models by assigning values to each unit. Once the physical property models were created, forward modeling of the conductivity model was performed in order to create synthetic surface data. Following the simulation of surface data, the methodology outlined in figure 5.1 can be applied.

5.1.1 Creating Geological Models

The synthetic modeling process began with the derivation of a geological model. In this step, the geometry was specified, including the size and complexity (ie: number of distinct geological units) of the model, as well as the depth of targets. In this research, 2D modeling was chosen so as to facilitate the computational aspects of the modeling (ie: run time), as well as the visualization of the results. Models were created with indices (ie: 1,

2, 3 etc.) to identify distinct rock types such that physical properties could then be easily applied in the following step.

The first example has been kept very simple: an attempt was made to recover a conductive cylinder in a resistive half space. Specifically, the cylinder is centered at 250m at a depth of 200m, and has a radius of 100m. Though this basic task was one easily accomplished by existing inversion methodologies, the goal was not only to recover the target, but to recover correct resistivity values with a well defined boundary between the two units. Geometrically, the model can was represented by the following two units:



Figure 5.2: Geological model with two units (red and blue)

5.1.2 Creating Geophysical Models

In order to provide a multidimensional downhole data-set, three physical properties (and thus three models) were defined for each unit: density, in grams per cubic centimeter (g/cm^3), magnetic susceptibility, in standard units (SI), and electrical resistivity, in Ohm-meters ($\Omega \cdot m$)⁶. These physical properties were chosen because they provide sufficient contrasts between the different units, and as such encompass some of the more commonly used methods for locating common geophysical targets.

To translate the geological units from figure 5.2 into geophysical models, the two units have been defined as a simplified mineralized zone within a sandstone background. The physical property values assigned to each unit were

⁶Physical property logs were actually collected as conductivities, in Siemens per meter (^S/m), then converted to resistivities, since $\rho = \frac{1}{\sigma}$

defined as normal distributions with a mean (μ) and standard deviation (σ) , where the mean values have been taken from Carmichael [1989] and the standard deviations were assigned based on comparison with real physical property logs⁷. Applying this to the first example, the following physical property values were assigned (see table 5.1 below).

	Dei	\mathbf{nsity}	Magnetic S	usceptibility	Electrica	l Resistivity
Unit	(g/cm^3)		(SI)		$(\Omega \cdot m)$	
	μ	σ	μ	σ	μ	σ
1	2.8	0.40	3.75×10^{-4}	7.5×10^{-5}	1.0×10^{3}	2.0×10^2
2	5.0	0.25	2.50×10^{-3}	5.0×10^{-4}	1.0×10^{1}	2.0×10^0

Table 5.1: Model I: Physical Property Values of Geological Units



Figure 5.3: Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density

⁷This method of assigning physical property values was chosen to better represent the variability present in a typical downhole log.

To simulate the high vertical resolution of downhole logging, the physical property models are defined on a fine mesh (figure 5.3 above). Before the models could be used for generating synthetic surface data, they were down-sampled onto a coarser mesh (figure 5.4) to facilitate the computational aspects of modeling.

As explained in chapter 4, simulating surface geophysical data can be accomplished through forward modeling of the physical property models. DC resistivity was the geophysical method of choice due to the superior depth of investigation as compared to potential field methods. The resistivity model was therefore forward modeled to simulate a DC resistivity survey using DCIPF2D (Uni [2011]).



Figure 5.4: True model

Specifically, two pole-dipole surveys (one in each direction) were simulated (see figure 5.5), each with 19 electrodes spaced 50m apart. The data were combined to increase the density of measurements and provide greater resolution at depth.



Figure 5.5: Pole-dipole survey. Top with pole current electrode (red) on left, bottom with pole current electrode (red) on right

Gaussian noise was added to the simulated data to represent instrument variability, and the associated standard deviations were assigned to the data as errors. For all data used, noise was also applied as a floor plus 5% of the datum. The noise floor was chosen such that the maximum noise level (on the smallest data) was approximately 30% of the datum value.



Figure 5.6: DC resistivity data (top) and associated percent errors (bottom)

5.1.3 Blind Inversion

Following the simulation of a DC resistivity survey, the data were inverted blind, with no prior information of the expected model. The inversion was run using DCIP2D (Uni [2011]), with the only defined parameters having been the mesh (the same mesh as was used for the forward modeling) and the data from the forward modeling. Since DCIP data can be susceptible to noise on the surface (near electrodes), it is common to add surface weighting to the inversion so that the top few layers of the model are encouraged to be smooth horizontally. This has the effect of reducing the amount of noise in the model resulting from the placement of the electrodes.





Figure 5.7: Blind inversion, with (bottom) and without (top) horizontal surface weighting

The inversion ran for 9 iterations and achieved the target misfit of 380 with a best fitting half-space of $515\Omega \cdot m$ used as a reference model. figure 5.7 above shows the resulting model with and without horizontal surface weighting, both with a color scale ranging from $1.0\Omega \cdot m$ to $1500\Omega \cdot m^{-8}$. Both recovered models are shown to isolate the impact made by the surface weighting.

The recovered model (figure 5.7a) clearly defines a body of some resistivity ranging from tens to hundreds of $\Omega \cdot m$ in magnitude, beginning at approximately 100m and extending down some few hundred meters. The width

⁸For consistency, when possible, models for each example will be displayed with the same colorscale, and always with padding cells removed

5.1. Methodology

of the body is somewhere in the range of 500m, centered at 250m, and the background resistivity value is approximately $550\Omega \cdot m$, with a slightly more resistive overburden.

This result was taken as a preliminary investigation into the region, and though the resulting model was not to be entirely trusted, it was used to inform the decision as to where to drill and log boreholes in the following step.

5.1.4 Simulating Physical Property Logs

Given this initial inversion, it was decided that extra information would be acquired via downhole physical property logging. In order to better define the magnitude and extent of the target, drill holes were placed at 210m, 250m, and 290m. This number of holes was chosen so as to allow for easy visualization of the results, while providing sufficient downhole information to successfully complete the procedure. Conductivity, magnetic susceptibility, and density were logged in each hole on 1m intervals from the surface to the bottom of the model, at 500m.

Since many downhole physical property measurements are relatively accurate in the near-field, downhole logs were simulated by extracting a column of data values from the fine resolution physical property models⁹, and applying a noise floor and percent error to represent instrument variability. To simplify the simulation, all boreholes were logged vertically (no dip or azimuth), and are of equal depth.

⁹Due to the large range in conductivity and magnetic susceptibility values, these physical properties were taken as the log of the model value. Similarly, since it is common practice to observe the anomalous density rather than the absolute density, anomalous density was logged by subtracting a background value from all measurements.

5.1. Methodology



Figure 5.8: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)

5.1.5 Classifying Downhole Data

Though it would be possible to directly apply the downhole physical property logs as a constraint for the inversion of surface data through the use of a reference model, it is not ideal. This is because of the extreme difference in resolution between downhole data and surface data, resulting in the potential for a highly variable reference model. Additionally, such a methodology requires some scheme of interpolation, which invariably requires the user to specify some level of information as to how to spread the borehole information out from the hole. This can be highly subjective and lead to difficulties, since borehole measurements are only sensitive to areas close to the hole.

Since it is common practice to collect multiple physical property logs simultaneously, statistical classification is an ideal candidate to simplify the high resolution information from the downhole logs and provide the inversion with coherent constraints. This is done using the Expectation-Maximization algorithm, as explained in section 3.3. The classification scheme takes as inputs the physical property logs (for as many holes and properties as are available) and returns as outputs a class ID for each datum, as well as an associated dictionary which translates each class into a mean value and standard deviation for each physical property. For the current example, the logs were classified for different numbers of rock types, ranging from a single unit up to six distinct rock types. Looking at a plot of the Akaike Information Criterion versus the number of clusters (figure 5.9), it was clear that the optimal number of clusters was two, as per section 3.3.7.1.



Figure 5.9: Akaike information criterion vs number of clusters

Given the classification results, this model has been defined in terms of two rock types, the first with a mean conductivity of $9.6\Omega \cdot m$ and the second of $857.0\Omega \cdot m$.



Figure 5.10: Classification results. Left to right: rock type vs depth & scatter plot

5.1.6 Creating Constraints

Once the data was classified, the results were applied to create constraints for geophysical inversion. In order to create constraints which are both reliable and applicable to the whole model, it has been assumed that all distinct rock types found in the model have been sufficiently sampled in the downhole logs, and therefore that the classification results are themselves representative of the true model.

As was previously mentioned, in order to minimize the amount of required user input, emphasis was put on the creation of upper and lower bounds files, which specify the minimum and maximum possible values for each cell of the recovered model. For consistency among constraints, a simple reference model was also generated rather than to rely on the best fitting half space determined by the inversion algorithm.

5.1.6.1 Generating Bounds

The initial upper and lower bounds files were kept fairly simple. For all areas surrounding the boreholes, bounds were assigned as the most extreme values possible (lower bound was the smallest possible value, upper bound was the largest possible value). The smallest and largest values, which will be referred to as ρ_{min} and ρ_{max} , were defined by the classification results. Since each class, or rock unit, was defined as a normal distribution, it was assumed that the majority of values belonging to each rock type were within one standard deviation (σ) of their mean value (μ). The smallest and largest values were then defined to be the minimum value of $\mu - \sigma$ and the maximum value of $\mu + \sigma$, respectively, for all rock units. For the first example, this equates to lower and upper bounds for resistivity of $6\Omega \cdot m$ (yellow) and $1240\Omega \cdot m$ (blue), respectively¹⁰.

¹⁰These values will be re-used later during the updating of constraints for cells which could not be classified.





Figure 5.11: Initial resistivity bounds constraints: upper bounds (top) and lower bounds (bottom)

With ρ_{min} and ρ_{max} set as the background resistivity values of the lower and upper bounds, respectively, it can be assumed that the majority of values in the recovered model would fall within these bounds, without concern of restricting the ability of the inversion to fit the data or add supported structure.

Additionally, since the values measured in the boreholes were known to be fairly accurate, the upper and lower bounds of cells in the path of the borehole were assigned a value of $\mu_i \pm \sigma_i$ (for upper and lower bounds, respectively), where *i* is the class ID of the cell as assigned during classification. In this way the inversion was entrusted to spread the borehole information out

through the model, without biasing it towards or away from the true model.

5.1.6.2 Generating a Reference Model

Similar to the bounds files, the reference model was created as a background value with the borehole values from classification overlayed. The background value was determined as the mean value of all measured borehole values ($433\Omega \cdot m$ for this example), and the borehole values were defined as μ_i , where *i* is the class ID of the cell as assigned during classification. Though this reference model was almost certainly wrong for most cells, it represents an approximate value within the range of the upper and lower bounds, and as was previously mentioned, the emphasis is on the bounds files rather than the reference model.



Figure 5.12: Reference resistivity model

5.1.7 Inversion With Constraints

In the second round of inversion, this time with constraints, a number of extra parameters were applied as compared with the initial blind inversion:

- The bounds files were applied as upper and lower bounds
- The classification derived reference model was used as an initial and reference model

- To minimize the effects of the reference model, α_s was set very small $(\sim 10^{-9})$ and
- USE_MREF was specified as FALSE so that the reference model was not used in the derivative terms of the objective function¹¹.

Given the previous inversion result as well as the classification results (figure 5.10), it was advisable to add in weighting files to mitigate excessively noisy structure on the surface that is typical of the inversion of DC data¹². Other than the aforementioned changes, the rest of the input file was the same.



Figure 5.13: Recovered model from inversion with borehole constraints

When compared with the initial inversion result presented in figure 5.7, this is a noticeable improvement in resolution. The target is now discernible as a somewhat circular body (in 2D cross-section) of approximately $15\Omega \cdot m$, ranging in depth from 100m to approximately 325m, with a lateral extent from approximately 100m to 400m. The background now has a more resistive value, with a blurred lateral band of higher resistivity (~ $1200\Omega \cdot m$) from approximately 150m to 250m.

¹¹In this way, the recovered model was only very slightly penalized for deviating from the reference model, allowing the emphasis to be placed on the bounds.

¹²The effects of the surface weighting are shown in figure 5.7

5.1.8 Updating Constraints, Re-inverting & Iterating

Though the incorporation of constraints greatly improved the resolution of the recovered model, further refinements were still possible. To accomplish this, an iterative process of classifying the resulting model and updating the upper and lower bounds files followed.

5.1.8.1 Classifying the Inversion Model & Updating the Bounds

The classification of the inversion model, much like the bounds and reference model creation, relied on the results from the initial classification of the downhole physical property logs. It was a three-tiered process that operated in the following manner:

- 1. Bin the recovered model
- 2. Overlay borehole values
- 3. Assign all unclassified cells wide bounds $(\rho_{min} \& \rho_{max})$

First, for each rock unit (i) in the dictionary, all cells in the recovered model with values which fell within C_1 standard deviations of the mean value of that rock unit were assigned a new upper and lower bound of $\mu_i \pm C_2 \sigma_i$ (figure 5.15b). Mathematically this is the following:

For each rock unit i, find all cells in recovered model m such that

$$(\mu_i - C_1 \sigma_i) \le m \le (\mu_i + C_1 \sigma_i) \tag{5.1}$$

and assign these cells upper and lower bounds of

$$\mu_i \pm C_2 \sigma_i \tag{5.2}$$

Thus if $C_1 > C_2$ the value of the bounds are truncated at known values determined from classification of the physical property logs¹³.

¹³ Typically, $1 \le C_1 \le 3$, and $C_2 = 1$.



Figure 5.14: 1D diagram illustrating the model classification procedure. a) Binning the recovered model into rock types based on statistical classification results. b) Resulting upper and lower bounds. c) Expansion coefficients. d) Resulting upper and lower bounds when b) is multiplied by c).

In order to allow for flexibility in the recovered model, the newly assigned upper and lower bounds were multiplied by an expansion coefficient, C_3^{14} (figure 5.15c), such that the assigned bounds expanded exponentially as the distance between the cell and the closest borehole increased. This coefficient was scaled such that the values of the resulting bounds ranged from the original assignment (from 5.2) to $2(C_3 - 1)$ times wider bounds. In this way, as the distance from a borehole increased, and the reliability of the model classification decreased, the bounds widened accordingly, allowing for a wider range of structures and values in the recovered model (figure 5.15d).

¹⁴Reasonable values for C_3 are $1.0 \le C_3 \le 2.0$



Figure 5.15: Model classification procedure. From top to bottom: previous recovered model, upper resistivity bounds without expansion coefficients, expansion coefficients, upper resistivity bounds with expansion coefficients.

Next, the bounds of the cells in the path of the boreholes were once more set to $\mu_i \pm \sigma_i$ (where *i* is the class assignment from the original classification). Since the downhole physical property logs were trusted to be accurate in close proximity to the holes, there was no reason to change the values assigned during the initial classification.

Finally, for all unclassified cells, the extreme bounds from the original bounds constraints, ρ_{min} and ρ_{max} , are applied. In this way, the model was left flexibility to alter the value of these unclassified cells until they fall within range of one of the classified units.

For the current example, the first iteration of classification was performed using

$$C_1 = 3.0$$

 $C_2 = 1.0$ (5.3)
 $C_3 = 1.5$

This implies that all values within three standard deviations of the mean value of one of the rock units will be binned as belonging to that unit, and the bounds will be reassigned as the mean value of that unit plus or minus one standard deviation (for upper and lower bounds, respectfully). Following this, these new bounds will be multiplied by a Gaussian weighting ranging from 1.0 at the borehole locations, out to 1.5 at the furthest point. As one can see in figure 5.16 below, most of the model was binned into one of the two rock units, with a buffer region between the two where the values were neither one nor the other. These cells were assigned the same values as were used for the upper and lower bounds backgrounds during the initial creation of constraints in section 5.1.6.





Figure 5.16: Updated upper (top) and lower (bottom) bounds, iteration 1

5.1.8.2 Re-inverting & Iterating

With the constraints updated such that the bounds for most cells have become tighter, the inversion was run once more, using the same parameters as the previous inversion, but with updated constraints. Applying these new bounds files as constraints for another round of inversion, the model below in figure 5.17 was recovered. Though it bears a close resemblance to the result from the initial round of constrained inversion (figure 5.13), there exists a more well defined boundary between the two units and the values of all cells are closer to the mean values defined in classification. Additionally, the background values have smoothed out to a more homogenous mean value of approximately $700\Omega \cdot m$, and the lateral band of resistive cells has somewhat dissipated.



Figure 5.17: Recovered model, iteration 1

The results from this new inversion were then classified just as were the previous results, and the bounds updated once more. In this way, the number of cells with wide bounds decreased with each iteration, and the inversion was able to converge to a model in which the value of each cell was within C_2 standard deviations of the mean value of one of the classified rock units. For the example being discussed, three iterations were performed. To speed up the process, at each iteration, C_3 was decreased so that the value of the bounds expanded less from the assigned value at each iteration. For the current example, the expansion coefficient (C_3) was decreased from 1.5 to 1.2 to 1.1, finally resulting in the following recovered model.





Figure 5.18: Recovered model, iteration 3 (top) compared to recovered model from blind inversion result (bottom)

If compared to the previous results, especially the initial blind inversion, this final recovered model has greatly improved the understanding of the exploration target, to the extent that it would be possible to define a boundary to within a couple of cells between a conductive body (comfortably estimated to be approximately $10\Omega \cdot m$) and a resistive background (of approximately $1000\Omega \cdot m$). Without the additional information provided by the downhole physical property logs, such accuracy would never have been possible, particularly at depth.

In the following subsections, three simple geological models are processed using the methodology just discussed.

5.2 Model I: Cylinder in a Half-space

5.2.1 The Model

Model I is a refinement on the example used in the previous section. The model is identical, as are the forward modeled surface data (for easy reference see figures 5.19, 5.20 and 5.21 below).



Figure 5.19: Model I: True model



Figure 5.20: Model I: DC resistivity data (top) and errors (bottom)

The difference in this example is that the boreholes will be taken in different locations to simulate a more realistic exploration project in which only one drill hole was able to pierce the target.



Figure 5.21: Model I: Blind Inversion with (bottom) and without (top) horizontal surface weighting

5.2.2 Downhole Physical Property Logs & Classification

Specifically, drill holes were placed at 100m, 250m, and 400m. Conductivity, magnetic susceptibility, and density were again logged in each hole on 1m intervals from the surface to the bottom of the model, at 1000m (including padding cells, not seen).



Figure 5.22: Model I: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)

The physical property logs (figure 5.22) were again statistically classified via the Expectation-Maximization algorithm for different numbers of rock types, ranging from a single unit up to six distinct rock types. Looking at a plot of the Akaike Information Criterion versus the number of clusters (figure 5.23) , the optimal number of clusters was still two, as per section 3.3.7.1.



Figure 5.23: Model I: Akaike information criterion vs number of clusters

Clustering all the downhole measurements into two distinct rock types, the following plot of rock type with depth produced was produced, with the associated scatter plot (figure 5.24):



Figure 5.24: Model I: Classification Results. Rock type vs depth (left) & scatter plot of physical property values (right)

Given the new classification results, Model I has been defined in terms of two rock types, the first with a mean resistivity of $9.4\Omega \cdot m$ and the second of $850.9\Omega \cdot m$. These values are very similar to the values from the previous example, and indicate that the statistical classification is not heavily impacted by the location of the drill holes, so long as the same rock types are still sampled sufficiently.

5.2.3 Creating & Applying Constraints

The classification results were applied as per section 5.1.6 to create constraints for the next round of inversion. When completed, the absolute lower and upper bounds for resistivity were $7.8\Omega \cdot m$ and $1021.1\Omega \cdot m$, respectively¹⁵. The reference model background was between these values, with a resistivity of $430.1\Omega \cdot m$. All three of these models can be seen below in figure 5.25.

 $^{^{15}{\}rm These}$ values will be re-used later during the updating of constraints for cells which could not be classified.



Figure 5.25: Model I: Initial constraints. From top to bottom: upper resistivity bounds, lower resistivity bounds, and reference model.

Applying these constraints, along with the surface weighting to smooth out the top 50m and a few other parameters discussed in section 5.1.7, the following model is recovered from the inversion of the DC resistivity data:



Figure 5.26: Model I: Recovered model from inversion with borehole constraints

Even with only a single borehole pierce point, the result presented above is a significant improvement compared to the original recovered model (figure 5.21). The model recovered with the help of the borehole constraints presents a much more finite body, beginning at approximately 100m and extending down in two lobes to a maximum depth of 350m. The width of the body has been greatly reduced, now ranging from 125m to 375m, with resistivities from $7 - 125\Omega \cdot m$. The background has become more homogenous thanks to the surface weighting, with an average value of around $850\Omega \cdot m$.

5.2.4 Updating Constraints and Final Models

In the next sections, the goal was to refine the result from the previous inversion. This was done iteratively using the results from classification as per the steps outlined in section 5.1.8, beginning with the model shown above in figure 5.26.
Iteration I

Classification of the model in figure 5.26 was again performed using an expansion factor C_3 of 1.5 (50% increase in bounds as the distance from boreholes increase) and a threshold of $\mu \pm 3\sigma$ ($C_1 = 3$, since the two units had such small variance). Visible in figure 5.27 below, most of the model was binned into one of the two rock units, though there exists a buffer region between the two where the values were neither one nor the other. These cells were assigned the same values as were used for the upper and lower bounds backgrounds (ρ_{min} and ρ_{max}) during the initial creation of constraints in section 5.2.3.



Figure 5.27: Model I: Updated bounds, iteration 1

These new bounds files were applied as constraints for another round of inversion, giving the model below in figure 5.28. Despite similarities to the result presented in figure 5.26, there exist subtle differences. Notably, the two lobes which extend down beneath the main body have been truncated, and the boundary between the central target and the surrounding background are now more well defined, with the values of all cells closer to the mean values defined in classification.



Figure 5.28: Model I: Recovered model, iteration 1

Iteration II

This time classification of the model in figure 5.28 was performed with a smaller expansion coefficient, $C_3 = 1.2$ (20% increase in bounds as the distance from boreholes increase). The same threshold was applied, and the bounds were updated to the following:



Figure 5.29: Model I: Updated bounds, iteration 2

During this iteration of classification less the buffer region between the two units is slightly larger, with a few additional patches of unclassified cells having been given the initial background bounds from figure 5.25. The classification had the most difficulty binning the edges of the conductive body in the center, however the majority of the background unit was correctly binned as such. The overall tighter bounds (20% expansion rather than 50%) encouraged the inversion to refine the model yet further, resulting in the following model:



5.2. Model I: Cylinder in a Half-space

Figure 5.30: Model I: Recovered model, iteration 2

Iteration III

Finally, in the third iteration of this procedure, the most recent inversion result was classified with an expansion coefficient of only $C_3 = 1.1$ (10% increase in bounds), and again the same threshold. The results are much the same as the last iteration for both the bounds (figure 5.31) and the inversion result (figure 5.32), and thus the iterative process was stopped.

In contrast to the initial blind inversion, this final recovered model has greatly improved the understanding of the exploration target. Though the result is not quite as definitive as the initial example shown, it is still possible to outline an elliptical body of approximately $10\Omega \cdot m$ sitting in a background or approximately $1000\Omega \cdot m$. Just as before, the addition of the information from downhole physical property logs has greatly improved the understanding of the model, even with only a single borehole pierce point.



5.2. Model I: Cylinder in a Half-space

Figure 5.31: Model I: Updated bounds, iteration III



Figure 5.32: Model I: Recovered model, iteration III (top), compared to blind inversion result (bottom)

These results will be further discussed in following chapters.

5.3 Model II: Vertical Contact with Resistive Overburden

5.3.1 The Model

This model introduces a new added level of difficulty: overburden. First, an attempt was made to recover a vertical contact between two units with a resistive overburden, and then in the next model the same scenario but with a conductive overburden. The two basement units extend from 50m below the surface to depth with a vertical contact in the center of the model, at 250m. On top of these, the top 50m of the model consists of a layer of resistive material. The difficulty in this example was recovering the correct location for the contacts, as well as the correct magnitude of resistivity values, despite the overburden. The geometry of the model is the following:



Figure 5.33: Models II & III: Geological model with three units (red, green and blue)

Translating this geological model into geophysical models, the overburden was interpreted to represent a layer of resistive soil, or till, sitting atop the contact between a sandstone unit (right) and a shale unit (left). The mean physical properties for these units are shown below in table 5.2 with associated standard deviations.

	Density		Magnetic Susceptibility		Electrical Resistivity	
Unit	(g/cm^3)		(SI)		$(\Omega \cdot m)$	
	μ	\sum	μ	\sum	μ	\sum
1	2.8	0.40	3.75×10^{-4}	$7.0 imes 10^{-5}$	1.0×10^3	$2.0 imes 10^2$
2	2.5	0.20	6.00×10^{-4}	1.0×10^{-4}	2.0×10^2	4.0×10^1
3	1.7	0.25	6.20×10^{-5}	2.0×10^{-5}	2.5×10^3	4.5×10^2

5.3. Model II: Vertical Contact with Resistive Overburden

Table 5.2: Model II: Physical Property Values of Geological Units

Applying these physical properties to the geological model in figure 5.33, the following three physical property models were created:



Figure 5.34: Model II: Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density

Once re-meshed, the true resistivity model which we attempted to recover was the following:



Figure 5.35: Model II: True model

5.3.2 DC Data & Blind Inversion

5.3.2.1 Data

Data was again forward modeled using DCIPF2D with a synthetically generated pole-dipole survey consisting of 19 electrodes space 50m apart (figure 5.5). For this data set the Gaussian noise was added with a floor of 0.0005 plus 5% of each datum. The resulting data and associated percent errors are shown below in figure 5.36.



Figure 5.36: Model II: DC resistivity data (top) and associated percent errors (bottom)

5.3.2.2 Inversion

The first inversion was run blind, with only the mesh and the data as inputs. After running for 9 iterations and achieving the target misfit of 380 with a best fitting half space of $420\Omega \cdot m$ used as a reference model, the inversion recovered the following models (one with and one without horizontal surface weighting):



Figure 5.37: Model II: Blind Inversion with (bottom) and without (top) horizontal surface weighting

These results clearly depict a resistive (in the range of $3000\Omega \cdot m$) overburden approximately 50m thick, with a large conductive body below it on the left. This body has a minimum resistivity of around $100\Omega \cdot m$ and begins at a depth of approximately 100m, extending down to at least 300m before fading to background values (in the range of $500\Omega \cdot m$). On the right there is a band of more resistive material that appears to be associated with the overburden, though it is smaller in magnitude.

5.3.3 Downhole Physical Property Logs & Classification

In order to further probe the conductive target and to better understand the resistive band on the right, three holes were drilled and logged: one at 60m, one at 185m, and the last at 340m. Again, the holes were logged for conductivity, magnetic susceptibility and density on 1m intervals ranging from the surface to the bottom of the model.



Figure 5.38: Model II: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)

Despite the initial inversion result, the logs indicated that the body on the left was actually a full unit, and gave no evidence of a resistive band on the right. To better explore this possibility, the downhole logs were classified using the Expectation-Maximization algorithm and applied to create constraints for a new round of inversion. The classification was performed for number of rock types, ranging from one to six, and the following Akaike Information Criterion vs number of clusters plot was produced, suggesting that there are three distinct rock types in the model.



Figure 5.39: Model II: Akaike information criterion vs number of clusters

Clustering the information from the physical property logs into three rock units, the following scatter plot and plot of rock type vs depth were produced:



Figure 5.40: Model II: Classification results. Rock type vs depth (left) & scatter plot of physical properties (right)

From the classification, it was determined that the resistivities of the three rock types are $1872.2\Omega \cdot m$, $868.7\Omega \cdot m$ and $190.9\Omega \cdot m$, respectively.

5.3.4 Creating & Applying Constraints

As per section 5.1.6, the results from classification were applied to create constraints consisting of upper and lower bounds as well as a reference model. The background lower and upper resistivity bounds were assigned as $125\Omega \cdot m$ and $2470\Omega \cdot m$, respectively, while the background reference model value was $977\Omega \cdot m$. The models, with borehole classification overlayed on these backgrounds, are visible below in figure 5.41.



Figure 5.41: Model II: Initial constraints. From top to bottom: upper resistivity bounds, lower resistivity bounds, and reference model.

Since the initial recovered model suggested a distinct layer of resistive material on the surface, which was further supported by the borehole logs, the surface weighting was altered to to smooth the top 50m and to allow a vertical discontinuity near the bottom of this layer. When this weighting alone is used, the following model is produced:



Figure 5.42: Model II: Recovered model from blind inversion with full surface weighting



Figure 5.43: Model II: Recovered model from inversion with borehole constraints and surface weighting

Combining this weighting with the aforementioned constraints and the parameters suggested in section 5.1.7, the model in figure 5.43, above was recovered.

Comparing this new result to the previous blind inversion, much of the character is the same, however some important features are different. As was suggested by the logs, the conductive body on the left continues to depth and has a far more homogenous structure than was previously expected, with an average resistivity in the range of $150\Omega \cdot m$. Additionally, the contact between the conductive unit and the overburden is much sharper, and nearly horizontal. The resistive body on the right has somewhat dissipated, while the overall resistivity of the right side of the model has increased from $\sim 500\Omega \cdot m$ to $\sim 800\Omega \cdot m$. Due to the surface weighting, there is much less noise in the top 50m, and the average resistivity of the overburden now sits at $\sim 2500\Omega \cdot m$.

5.3.5 Updating Constraints and Final Models

Beginning with the previous inversion result and iterating, an attempt was made to further increase the resolution of the recovered model.

Iteration I

As was previously done, the first iteration of classification used an expansion coefficient of $C_3 = 1.5$ (50% increase of bounds) with a threshold of $\mu \pm 2\sigma$ ($C_1 = 2$). Most of the model was classified, however a vertical band in the center of the model was not able to be binned as any of the three rock units, and thus was assigned the default background resistivity bounds ($\rho_{min} \& \rho_{max}$) from the initial bounds constraints (figure 5.41).



Figure 5.44: Model II: Updated bounds, iteration 1

Using these updated constraints for another iteration of inversion, the following model was recovered:



Figure 5.45: Model II: Recovered model, iteration 1

This recovered model continues the trend from the past inversion. The conductive unit on the left has become yet more homogenous, the average resistivity increasing further to approximately $200\Omega \cdot m$, while the resistive body on the right has dissipated yet more, and the resistivity of the background unit on the right has increased again: now $\sim 1000\Omega \cdot m$. Additionally, the contact between the overburden and the bottom two units has become very clean, and the contact between the two bottom units has moved more to the center of the model and become more vertical.

Iteration II

In this iteration of classification an expansion coefficient of $C_3 = 1.2$ was used (20% increase of bounds), again with the same threshold. The entire model was successfully binned in one of the three units this time, defining a very vertical boundary in the center of the model between the two basement units.



Figure 5.46: Model II: Updated bounds, iteration 2

With the tighter bounds applied (using only 20% expansion), the boundary between the two basement units has become quite sharp. The rest of the model is more or less unchanged, though the resistive body on the right has further dissipated.



5.3. Model II: Vertical Contact with Resistive Overburden

Figure 5.47: Model II: Recovered model, iteration 2

Iteration III

Finally, reducing the expansion coefficient to $C_3 = 1.1$, the bounds become very tight, though they bear the same overall geometry as the previous iteration.



Figure 5.48: Model II: Updated bounds, iteration 3



Applying these bounds to the third iteration, the final model was recovered:

Figure 5.49: Model II: Recovered model, Iteration 3 (top) compared to blind inversion result (bottom)

This final model maintains the sharp boundary between the two basement units, but also extends the conductivity of the left unit and the resistivity of the right unit to depth more than in previous models. Thus through iterating and updating the constraints, the contacts between all three units have been successfully resolved and the correct resistivity values determined. These results will be further discussed in a later chapter.

5.4 Model III: Vertical Contact with Conductive Overburden

5.4.1 The Model

Geometrically, the final model is the same as the last: two basement units extending from 50m below the surface to depth, with an overburden on top of them (see figure 5.33). The difference now is that rather than a resistive overburden, it is now a conductive overburden. The two basement units remain the same. Again, the goal was to recover the contacts between the units to as high resolution as possible, while correctly estimating the resistivity of each unit.

As opposed to the previous model, in which the overburden consisted of sand/till, in this model the two basement units are overlayed by a conductive layer of caliche. As such, physical property values were assigned based on typical values for such rock types Carmichael [1989] (see table 5.3 below).

	Density		Magnetic Susceptibility		Electrical Resistivity	
Unit	(g/cm^3)		(SI)		$\Omega \cdot m$	
	μ	\sum	μ	\sum	μ	\sum
1	2.8	0.40	3.75×10^{-4}	7.0×10^{-5}	1.0×10^3	2.0×10^2
2	2.5	0.20	6.00×10^{-4}	1.0×10^{-4}	2.0×10^2	4.0×10^1
3	1.9	0.25	9.00×10^{-4}	1.0×10^{-4}	4.0×10^2	$8.0 imes 10^1$

Table 5.3: Model III: Physical Property Values of Geological Units

When these values were applied, the following three physical property models were produced:



Figure 5.50: Model III: Physical property models on fine mesh. From top to bottom, resistivity, magnetic susceptibility, and density



Finally, the down-sampled true resistivity model is the following:

Figure 5.51: Model III: True model

5.4.2 DC Data & Blind Inversion

5.4.2.1 Data

Just as with the past two models, a DC resistivity survey was synthetically collected, again using a mirrored pole-dipole configuration with 19 electrodes spaced 50m apart (see figure 5.5). For this data set, the Gaussian noise was added with a noise floor of 0.0005 and 5% of each datum. The resulting data and percent error are presented in pseudo-section in the following figure:



Figure 5.52: Model III: DC resistivity data (top) and associated percent errors (bottom)

5.4.2.2 Inversion

A first pass blind inversion of this data was run with no prior information or constraints; the only input parameters defined were the mesh and the data from figure 5.52.

The inversion ran for 8 iterations before achieving the target misfit of 380 with a best fitting half space of $335\Omega \cdot m$ used as a reference model. figure 5.53 below shows the recovered model, with and without horizontal surface weighting applied, with padding cells removed and a color scale ranging from $150\Omega \cdot m$ to $2500\Omega \cdot m$.



Figure 5.53: Model III: Recovered model from blind inversion with (bottom) and without (top) horizontal surface weighting

From the resulting models, it is apparent that there exists a conductive body on the left with a resistive body on the right, though much else is not clear. The surface is very noisy, with extreme values in resistivity placed directly under the electrodes. The conductive body on the left extends from approximately 75m down to approximately 250m, with a lateral extent from -250mto 200m, with an average resistivity in the range of $150\Omega \cdot m$. The resistive body similarly begins at around 75m, though it is thinner and disappears by 150m depth. Laterally, the resistive body extends from around 300m, to the far right edge of the model at 750m (and on into the padding cells). The rest of the model depicts a smooth background with a resistivity in the range of $300 - 400\Omega \cdot m$.

5.4.3 Downhole Physical Property Logs & Classification

Based on the initial results, three holes were drilled to further explore the two recovered exploration targets (one resistive, one conductive). Since the conductive target was of more interest, two drill holes were designed to pierce this body, one at 75m and another at 175m, while one exploratory hole attempted to pierce the resistive body at 325m. From this information it was hoped that a better understanding of the interaction of these two bodies would be achieved. As with Model I, all three holes were logged for conductivity, magnetic susceptibility and density on a 1m interval extending from the surface to the bottom of the model.



Figure 5.54: Model III: Physical property logs. Properties listed above, units listed below, for three holes (listed at top)

From the physical property logs, it appeared that the bodies extend to depth, and thus it was decided that a more detailed inversion was required. To accomplish this constraints were required, and therefore the logs were classified using the Expectation-Maximization algorithm. The logs were classified for different numbers of rock types, ranging from a single unit to six distinct units. Looking at the Akaike plot below in figure 5.55, it was decided that there exists three distinct rock units in the recovered model.



Figure 5.55: Model III: Akaike information criterion vs number of classes

Clustering the physical property logs into three rock types, the following plot of rock type with depth was recovered, with the associated scatter plot for the three physical properties (figure 5.56):



Figure 5.56: Model III: Classification results. Rock type vs depth (left) & scatter plot (right)

Given the classification results, Model III has three rock types, the first with a mean resistivity of $376.9\Omega \cdot m$, the second a mean of $971.2\Omega \cdot m$, and the third of $194.2\Omega \cdot m$.

5.4.4 Creating & Applying Constraints

The classification results were applied as per section 5.1.6 to create constraints for the next round of inversion. When completed, the absolute lower and upper resistivity bounds were $130\Omega \cdot m$, $1450\Omega \cdot m$, respectively. The reference model background was between these values, with a resistivity of $514\Omega \cdot m$. All three of these models can be seen below in figure 5.57



Figure 5.57: Model III: Initial constraints. From top to bottom: upper resistivity bounds, lower resistivity bounds, and reference model.

In addition to these constraints the surface weighting was altered to smooth out the top 100m as well as allow the inversion to have a sharper boundary at the tops of the two bodies. This weighting was motivated by the results of classification which seem to indicate a sharp horizontal contact between the two bottom units and the overburden. Applying only this new weighting to the inversion, the following model is recovered:



Figure 5.58: Model III: Recovered model form blind inversion with full surface weighting

Combining all of these constraints and re-running the inversion with the parameters discussed in section 5.1.7, the following model was recovered:



Figure 5.59: Model III: Recovered model from inversion with borehole constraints

The latest inversion result presents a stark contrast to the initial blind inversion (figure 5.7) with two finite bodies in the near surface. As was apparent in the downhole logs, both bodies, particularly the conductive one, extend to depth, comprising not simply two bodies, but rather two distinct rock units with resistivities in the range of $150\Omega \cdot m$ and $700\Omega \cdot m$, respectively. From looking at this model it appears that the contact between the two dips at approximately 75^o , and that there still exists a resistive body ($\sim 1200\Omega \cdot m$) in the top right of the model. The surface appears significantly less noisy, with an average resistivity of approximately $400\Omega \cdot m$.

5.4.5 Updating Constraints and Final Models

In the next sections, the goal was to refine the result from the previous inversion. This was done iteratively using the results from classification, as per the steps outlined in section 5.1.8, beginning with the model shown in figure 5.59.

Iteration I

Classification of the model in figure 5.59 was performed as per section 5.1.8 with an expansion coefficient of $C_3 = 1.5$ (50% increase in bounds as the

distance from boreholes increases) and a threshold of $\mu \pm 2\sigma$ ($C_1 = 2$). Most of the model was easily binned into one of the three rock types, with a region between the two bottom units bearing the same rock type as the overburden. The weighting used to expand the bounds away from the holes is clearly visible, giving the model room to adjust the model as necessary in the next iteration.



Figure 5.60: Model III: Updated bounds, iteration 1

With the bounds files updated (figure 5.60 above), another round of inversion was completed, again with the same parameters as the last. The newly recovered model looks very similar to the last. The only differences being that the resistive body on the right appears to have somewhat dissipated and the contact between the two basement units is at a sharper angle.


5.4. Model III: Vertical Contact with Conductive Overburden

Figure 5.61: Model III: Recovered model, iteration 2

Iterations II & III

As per the suggested procedure, classification was again performed, first with a smaller expansion coefficient of $C_3 = 1.2$ (20% increase in bounds), and then $C_3 = 1.1$ for iteration III. Again, most cells were classified as being one of the three rock types in the dictionary, with a band running down the middle of the model classified as the same rock type as the overburden. Updating the bounds files, the following models were produced:



Figure 5.62: Model III: Updated bounds, iteration 2



Figure 5.63: Model III: Updated bounds, iteration 3

Since the geometry of the bounding models did not significantly changed, the geometry of the recovered models don't change very much either. The contact continues to become more vertical and the resistivity of the body on the right continues to disperse, while the unit on the left becomes less conductive. The region between the two appears as a buffer zone of values similar in resistivity to the overburden. Comparing figures 5.64 & 5.65, the recovered models from both iterations II and III are very similar.



5.4. Model III: Vertical Contact with Conductive Overburden

Figure 5.64: Model III: Recovered model, iteration 2



Figure 5.65: Model III: Recovered model, iteration 3

Iteration IV

When the expansion coefficient was finally dropped to $C_3 = 1.0$ (no expansion) the bounds files became very sharp. The contact between the two units is nearly vertical, while the contact between the overburden and the bottom two is nearly horizontal. There still exists a vertical band of cells binned as belonging to the same rock type as the overburden, however it is much narrower than was initially recovered.



Figure 5.66: Model III: Updated bounds, iteration 4

This final recovered model is a further progression of the other iterations: the contact is now practically vertical, and the resistive body at the top right is much smaller. The resistivity of the right unit is approximately $900\Omega \cdot m$, while the resistivity of the right unit has increased to approximately $200\Omega \cdot m$. The surface, as well as a thin vertical band separating the two basement units, both have a resistivity in the range of $400\Omega \cdot m$, and aside from a few anomalies, each of these units is fairly homogenous.



Figure 5.67: Model III: Recovered model, iteration 4 (top), compared to blind inversion result (bottom)

Comparing this result with both the initial blind inversion, as well as the first constrained inversion, the resolution has been greatly increased. Again, as with Models I & II, the addition of downhole physical property logs, paired with the iterative methodology suggested, has significantly improved the understanding of the model.

Chapter 6

Discussion

6.1 Summary of Results

Though the models used as examples in this thesis were relatively simple, they serve as a proof of concept that the suggested methodology can have dramatic benefits over existing unconstrained inversion. Below a summary of the results is presented, with a brief discussion on the benefits and challenges of the applied methodology.

6.1.1 Model I: Cylinder in a Halfspace

Model I was used as a first attempt due to the simplicity of the model: only two units, both with relatively low variability in physical properties (see figure 5.8). The goal with this model was to test the ability of the iterative scheme in refining the model, rather than to challenge the classification.

The initial blind inversion (figure 6.1, middle) was able to detect a body at the correct depth, however both the magnitude of the resistivity as well as the size of the body are poorly estimated. In addition to this, the magnitude of the background resistivity was approximately half of the true value, at only $550\Omega \cdot m$. Though arguments can be made for applying further constraints to the initial inversion rather than running it with all default values, with no prior knowledge of the expected geology, there is no basis for imposing constraints at this stage.

Once motivated by the initial recovered model, the added information from the borehole logs is extremely effective at refining the resolution of the target body (figure 6.1, bottom). The final iterations are an attempt to further increase the resolution, particularly at the boundary between background and target. Since the borehole information provides support for the target in the vertical direction, the vertical extent of the target is correctly recovered, with sharp boundaries on the top and bottom. However since the inversion is required to estimate the horizontal extent of the target away from holes, the horizontal extent of the target is slightly misrepresented. Part of this can be attributed to the fact both the borehole logs and the mesh are rectilinear while the target is circular; thus the sides of the target are less well approximated. Despite this, the final result is extremely close to the true model, with only $\sim 5 - 10\%$ of the target being poorly recovered (see figure 6.1 below).



6.1. Summary of Results

Figure 6.1: Model I: Results. From top to bottom: True model, model recovered from blind inversion, model recovered from suggested methodology. Outline of true lithological boundaries visible in black

6.1. Summary of Results

The results just presented bear more weight when one considers that the depth of investigation in these models is only the first couple hundred meters. This is due to the survey configuration (Oldenburg and Li [1999]) as well as to the resistivity values of the model. As a rule of thumb, the maximum depth of investigation for a pole-dipole array is taken to be between one third and one half of the maximum electrode spacing 1000m, which in the case of this survey (figure 5.5) is between 333m and 500m. This approximation should be taken with a grain of salt however, since other factors will play into the true depth of investigation. Additionally, it should be noted that this does not imply quality data down to a given depth, but rather decreasing resolution with depth.

In addition to this, difficulties can arise from a very resistive or a very conductive overburden. For very resistive overburden, it can be difficult to get sufficient current into the ground, resulting in a low signal to noise ratio. At the other end of the spectrum, if the overburden is too conductive, all of the current put into the ground can be channeled along the surface, resulting in little to no signal at depth.



Figure 6.2: Model I: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5

Average sensitivities have been calculated for the surface data. These values are normalized so that the sensitivity ranges from zero to unity, with the most sensitive cells in the model being at the surface (near the measurement locations). The contour in figure 6.2 above shows the depth at which the surface data is only 50% sensitive to structure.

6.1.2 Model II: Vertical Contact with Resistive Overburden

Model II was an investigation into the ability of constraints to assist in the recovery of a geological contact masked by an overburden. More variability was added to the model, and structure continued to depth. In this case the initial blind inversion was able to fairly accurately recover the geometry of the overburden, however the magnitude (off the color-scale in figure 6.3, middle) is far too large. Additionally, the bodies recovered beneath the overburden are quite misleading. Rather than a vertical contact between two units, it appears that there is a finite, oblong conductive body sitting in a fairly uniform background of approximately $400\Omega \cdot m$.

After surface weighting and borehole constraints are applied, the model changes significantly. By the end of the iterative procedure, a very different model is presented (figure 6.3, bottom). Both the vertical and horizontal contacts are accurately recovered, and each of the two basement units are much more homogenous than those recovered in the initial inversion.





Figure 6.3: Model II: Results. From top to bottom: True model, model recovered from blind inversion, model recovered from suggested methodology. Outline of true lithological boundaries visible in black

Though some artifacts exist, such as the more resistive region in the top left corner of the bottom right unit and the two conductive lobes in the bottom left unit, the primary features of the final recovered model are the three distinct units, with the vertical and horizontal contacts. Furthermore, these artifacts appear to be supported by the data, since they appear in all of the inversions, both with and without constraints. It is possible that they were introduced by a combination of factors, including the generation of models using normal distributions of values, as well as the random Gaussian noise applied to the collected surface data.



Figure 6.4: Model II: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5

Again, when the depth of investigation (figure 6.4, above) is considered, these results are fairly impressive. Given that the surface data quickly loses reliability after the first 200m, it is clear that the borehole constraints have added support for structure at depth.

6.1.3 Model III: Vertical Contact with Conductive Overburden

The third and final model was used as an extension of Model II. Model III investigated the effect of a conductive overburden over a vertical contact. This was chosen due to the difficulty that an inversion can have in recovering bodies beneath a conductive layer. In previous modeling attempts which contained conductive units on the surface overlaying finite conductive bodies at depth, virtually no signal was found in the DC resistivity data from the targets at depth. As a simplification of such a scenario, Model III (figure 6.5, top) was created.

Similar to Model II, the initial blind inversion of Model III recovers two finite oblong bodies in the shallow subsurface, sitting in a background of approximately $400\Omega \cdot m$. Contrary to the previous model however, the overburden is not accurately represented. Rather than a single homogenous band at the surface, the inversion has recovered a noisy mix of extreme values, both high and low resistivities. It is also worth noting that the recovered magnitude of the resistive body on the right is near $2000\Omega \cdot m$, nearly twice the highest resistivity found in the model¹⁶.

¹⁶Highest mean value of any of the three units is only $1000\Omega \cdot m$ (see table 5.2)



Figure 6.5: Model III: Results. From top to bottom: True model, model recovered form blind inversion, model recovered from suggested methodology. Outline of true lithological boundaries visible in black

X (m)

6.1. Summary of Results

By applying the borehole constraints, paired with the suggested iterative methodology, the final recovered model is very similar to the true model (see figure 6.5, bottom). Due to the surface weighting (motivated by the results of the physical property logs), the overburden is a much more distinct layer, with a fairly homogenous resistivity of $400\Omega \cdot m$. The two basement units are also much more apparent, the constraints having extended them to depth. Save for the resistive body in the top left of the bottom right unit¹⁷, the final recovered model is very accurate throughout, with physical property magnitudes and contact locations correctly estimated.



Figure 6.6: Model III: Depth of investigation. Contour at depth at which sensitivity has been reduced to 0.5

Once more, looking at figure 6.6 above, it is clear that the surface data does not provide reliable information below approximately 200m. Given this, one can conclude that the borehole constraints have added important and reliable information to the inversion, particularly at depth.

¹⁷Similar to in Model II, this artifact exists throughout all inversion, suggesting that it is supported by the data.

6.2 False Classification

One of the difficulties in applying classification to the discrimination of distinct rock units arises when the contrast between all values is very small. Under such circumstances, the classification of the recovered model (section 5.1.8) can falsely classify regions of the model surrounding the interface between two distinct units. This is particularly evident in the last two example models.

In Model II, the mean values of the bottom right, bottom left, and overburden units are $1000\Omega \cdot m$, $200\Omega \cdot m$, and $2500\Omega \cdot m$, respectfully. Since the value of the bottom right unit falls in between the two other units, the classification scheme has difficulty correctly binning the cells lying on the interface between the bottom left unit and the overburden (see figure 6.3, bottom). In Model II, this was partially remedied by the coincidental application of surface weighting to enhance the sharpness of this interface.

Model III suffers from the same problem, however the contrast has switched. With the mean resistivities of the bottom right unit, bottom left unit, and overburden at $1000\Omega \cdot m$, $200\Omega \cdot m$, and $400\Omega \cdot m$, respectively, it is now the vertical contact which suffers from this challenge in classification. In this case, no weighting was applied on this interface since there was no basis for forcefully imposing a horizontal discontinuity at a specific location in the model. The result is the narrow vertical band of yellow between the two basement units in the bottom image of figure 6.5. Looking back at section 5.4.5, it is clear that the methodology was able to reduce the impact of this effect through iterative classification, however due to this fundamental difficulty it is unable to completely resolve it.

One can imagine that as the model becomes more complex, and more distinct rock units are introduced, this will become more of an issue. One way to address this is to narrow the range of values to be binned as a given unit. Looking back at section 5.1.8, this would be accomplished by using a smaller value for C_1 such that only values very close to the mean value of a given unit are binned as belonging to it. This would have the effect of binning much less of the recovered model at each iteration, and accordingly would likely require many more iterations to converge to a final recovered model.

6.3 Choosing Parameters

One of the main motivations behind the derivation of the methodology used in this thesis was to remove bias from the user. It was stated that this would be accomplished by removing as much of the required user input as possible, and it may therefore seem counter-intuitive to impose parameters. The important difference is that the three parameters, C_1 , C_2 , and C_3 , act indirectly on the inversion. This is because the parameters set thresholds for the binning scheme presented in 5.1.8.1, which itself is statistically based. Additionally, the results of the binning scheme are the bounds constraints, which are applied via an iterative procedure, allowing flexibility for the methodology to remove as much of the user bias as possible.

Furthermore, the choice of the three parameters can be motivated by an analysis of the variability in the borehole logs, derived from the statistical classification. If the physical properties of the rock units are similar and highly variable, such that their signals overlap, it is desirable to apply smaller values of $C_1 (\leq 2)$ so that the risk of false classification is reduced. C_3 can also be made larger to allow for wider bounds far from boreholes, allowing greater flexibility to the inversion to recover from a potential false classification. This will likely result in a slower convergence and more iterations, but can produce more reliable results if units are highly variable. C_2 does not need to be changed, and can generally be kept as 1.

6.4 Real Data

Clearly the proposed methodology has had some success when applied to synthetic data, however the real test is to see how well it performs when applied to real geophysical data sets. The original intent of this thesis was to show both the synthetic examples as well as a real example. Due to logistical issues in acquiring surface geophysical data, this was never possible, however downhole physical property logs were supplied by DGI Geoscience Ltd.

Since no surface geophysical data was available at the time of writing, it was not possible to test the iterative inversion procedure. Despite a lack of surface data, statistical classification of the supplied downhole physical property logs was performed, with reasonable success. While the procedures and algorithms for classifying real data are identical to those applied for the synthetic data, working with real physical property logs introduces a few extra difficulties that must be taken into account.

6.4.1 Extra Considerations

For starters, selecting which physical properties to classify together is a non trivial task. Certain physical properties complement each other better than others, while some simply contain redundant information (ie: resistivity logs of different spacings). Methods exist to determine which logs contain the most information (S.E. MacMahon [2002]), though this was not applied in this thesis. In addition to this, many downhole logs are measured on multiple runs to ensure quality data. Picking the best run for a given log is not always straight-forward. Determining what is true signal, versus what is poor data can be a difficult and very subjective task.

Finally, the distribution of physical properties in the subsurface was modeled as being Gaussian, however this is a simple approximation. Whether or not it is a valid assumption, it is at best a simple approximation of a number of rather complicated geological processes. As such, when classifying real downhole physical property logs using the Expectation-Maximization method (which assumes the model is a set of multidimensional Gaussian distributions), there is often a much higher degree of striping in the classification results. This is to say that the resulting lithological model has many more thin, alternating layers of different rock types. To a certain degree, this simply represents the complexity of the true downhole structure, however below a certain vertical resolution, bands of a given rock type become nothing more than statistical noise. Since the downhole classification results will be applied as constraints for a much coarser model, isolated thin bands of a given rock type will be averaged out in favor of the bulk value.

6.4.2 Sample Classification

Here a brief summary of the results from a test run of classification of the real data for a single borehole are presented. Taking full advantage of the numerous physical properties logged, 8 physical property logs (Far density, P-wave slowness, magnetic susceptibility, temperature, neutron, gamma-gamma, fluid resistivity, and 8 inch normal resistivity) were used for statistical classification (figure below).



Figure 6.7: Sample of real physical property logs used in classification

These logs were classified using the EM method as described in section 3.3, for a model with the number of rock types ranging from 1 to 12. Plotting the Akaike Information Criterion vs the number of clusters produces the following curve:



Figure 6.8: Akaike Information Criterion vs. number of clusters for sample of real downhole data

Choosing to display the model found for 7 distinct rock types (notice the slight corner in the curve at 7 clusters), the following results are produced:



Figure 6.9: Results from EM classification of sample of real downhole data. Left, scatter plot of values for magnetic susceptibility vs. P-wave slowness vs. 8" normal resistivity. Right, plot of rock type vs depth.

Despite the extra striping (as discussed above), it is clear that there are 5 main units, with a couple of extra units coming in as bands. Plotting the recovered mean values (red) and standard deviations (green) versus the true data (blue) for each of the logs and recovered rock types gives the following:



Figure 6.10: Classification results for sample of real downhole data

Given the success of the algorithm in classifying real downhole physical property logs into distinct rock types, application of the full iterative inversion methodology would be a very interesting next step for this research, if and when the surface data should become available.

Chapter 7

Conclusions

This thesis undertook the challenge of developing a methodology whereby downhole physical property logs were leveraged to constrain geophysical inversion of surface data with minimal required user input. The resulting procedure employes an iterative scheme to refine the inversion model by updating the constraints via repeated re-classification of the recovered model. The constraints were applied mainly as upper and lower bounds files, with the emphasis removed from any reference model by setting α_s very small and removing the m_{ref} from the derivative terms in the objective function.

The suggested methodology was demonstrated on three synthetic geological models: a cylinder in a halfspace and a vertical contact with both a resistive and conductive overburden. The results of this research were compared to blind inversions which were run without any constraints. From these results, it can be concluded that the procedure is capable of increasing the resolution of recovered models in three main aspects:

- 1. The interfaces between units are sharper and more accurate
- 2. The physical property values of recovered bodies/units are both more homogenous and more accurate
- 3. Structures at depth are much more likely to be accurately recovered, despite limited sensitivity of the surface data

Given the success, it is be suggested that further research continue to explore each of the elements involved in this thesis. In particular, adapting the methodology to better handle the difficulties arising from false classification is important. Attempting this procedure with more complex models would be a worthwhile next step before complicating matters with further alterations. Many other aspects of this research could also benefit from more intensive scrutiny. For example, little effort has been put into determining which combinations of physical property logs are best for classification. Work on this has been done by using statistical methods to determine the relative importance of each log (S.E. MacMahon [2002]), and this would be a logical addition to the methodology.

To remove all user bias from the iterative scheme, a method of selecting the three parameters C_1 , C_2 , and C_3 directly from the statistical classification could be developed. Optimizing the combination of these parameters would greatly improve the iterative scheme and would greatly simplify the application of the procedure.

Furthermore, the Expectation-Maximization method was chosen as the classification method of choice due to its inherent flexibility and success in fitting the data during early experimentation, however it is possible that other methods would be better suited to this application, such as Self Organizing Maps (Kohonen [1990], Fraser and Hodgkinson [2009]).

For simplicity, this thesis applied only 2D inversion to the problem. To expand the investigation to 3D would be an interesting endeavor, and would almost certainly introduce further questions. Another possible application worth exploring is the merit of this methodology when combined with cooperative inversions. Since physical property classification inherently creates constraints which are consistent across physical properties, the method could be successfully applied to constraining cooperative inversions.

Bibliography

- Borehole geophysical methods. URL http://www.epa.gov/nerlesd1/cmb/ GeophysicsWebsite/pages/reference/methods/index.htm.
- Borehole geophysics and petrophysics. URL http://cgc.rncan.gc.ca/ borehole/site_e.php.
- Introduction to borehole geophysics. URL http://ny.water.usgs.gov/ projects/bgag/intro.text.html.
- Robert S. Carmichael. Practical Handbook of Physical Properties of Rocks and Minerals. CRC Press, 1989.
- Yihua Chen and Maya R. Gupta. Em-demystified: An expectationmaximization tutorial. Technical report, Department of Electrical Engineering, University of Washington, 2010.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38, 1977.
- Chuong B. Do and Serafim Batzoglou. What is the expectation maximization algorithm. *Nature Biotechnology*, 26(8):897–899, August 2008.
- Chris Fraley and Adrian E. Raftery. How many clusters? which clustering method? answers via model-based cluster analysis. *The Computer Journal*, 41:578–588, 1998.
- Stephen J. Fraser and Jane H. Hodgkinson. An investigation using sirosom for the analysis of quest stream-sediment and lake-sediment geochemical data. Technical report, Geoscience BC, September 2009.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The Elements of Statistical Learning. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.

- Thomas Hermans, David Caterina, Martin Roland, Andreas Kemna, Tanguy Robert, and Frederic Nguyen. How to incorporate prior information in geophysical inverse problems: deterministic and geostatistical approaches. In EarthDoc - Near Surface 2011 - 17th European Meeting of Environmental and Engineering Geophysics, 2011.
- Edward H. Isaaks and R. Mohan Srivastava. *Applied Geostatistics*. Oxford University Press, 1989.
- A.G Journel and C.J. Huijbregts. *Mining Geostatistics*. Academic Press Inc, 1978.
- P.G. Killeen. Borehole geophysics: Exploring the third dimension. In A.G. Gubins, editor, Proceedings of Exploration 97: Fourth Decennial International Conference on Mineral Exploration, pages 31–42, 1997.
- Teuvo Kohonen. The self-organizing map. *Proceedings of the IEEE*, 78: 1464–1477, 1990.
- Peter LeLievre. Integrating Geologic and Geophysical Data through Advanced Constrained Inversion. PhD thesis, University of British Columbia, 2009.
- Susanne MacMahon, Gary Hodgkinson, Dirk Kassenaar, and Bill Morris. Multiwell analysis of downhole physical rock properties of kimberlite: Guacho kue, northwest territories. In Proceedings of the 8th International KEGS/MGLS Symposium on Logging for Minerals and Geotechnical Applications, Toronto, August 2002.
- Georges Matheron. Principles of geostatistics. *Economic Geology*, 58(8): 1246–1266, December 1963.
- Geoffrey J. McLachlan and T. Krishnan. The EM algorithm and Extensions. Wiley, New York, 2 edition, 2008. ISBN 0471123587.
- Geoffrey J Mclachlan and David Peel. *Finite Mixture Models*. Wiley, New York, 2000.
- D.W. Oldenburg and Y. Li. Estimating depth of investigation in dc resistivity and ip surveys. *Geophysics*, 64(2):403–416, 1999.
- D.W. Oldenburg and Y. Li. Inversion for applied geophysics: A tutorial. Society of Exploration Geophysics Investigations in Geophysics, 13:89–150, 2005.

- D.W. Oldenburg and D.A. Pratt. Geophysical inversion for mineral exploration: a decade of progress in theory and practice. In B. Milkereit, editor, *Proceedings of Exploration 07: Fifth Decennial International Conference* on Mineral Exploration, pages 61–95, 2007.
- G. R. Pickett. Applications for borehole geophysics in geophysical exploration. *Geophysics*, 35(1):81-92, 1970. doi: 10.1190/1.1440083. URL http://link.aip.org/link/?GPY/35/81/1.
- B. M. P.K. Fullagar, G.A. Pears. Towards geologically realistic inversion. In Proceedings of Exploration 07: Fifth Decennial International Conference on Mineral Exploration, 2007.
- B. McMonnies P.K. Fullagar, G.A. Pears. Constrained inversion of geologic surfaces: Pushing the boundaries. *The Leading Edge*, January:98–105, 2008.
- Tapio Schneider. Analysis of incomplete climate data: Estimation of mean values and covariance matrices and imputation of missing values. *Journal of Climate*, pages 853–871, 2001.
- D.J. Kassenaar W.A. Morris S.E. MacMahon, G. Hodgkinson. Statistical analysis of downhole physical property measurements: Classification and predictive analysis. In *Proceedings of the 8th International KEGS BGLS* Symposium on Logging for Minerals and Geotechnical Applications, 2002.
- DCIP2D Manual. University of British Columbia Geophysical Inversion Facility, 5.0 edition, July 2011.
- Juha Vesanto, Johan Himberg, Esa Alhoniemi, and Juha Parhankangas. Som toolbox for matlab 5. Technical Report A57, Helsinki University of Technology, April 2000. URL http://www.cis.hut.fi/projects/ somtoolbox/.
- Peng Wang and Kai Yang. A geostatistical inversion technique constrained by well-log, cross-hole and surface seismic data based on visim. SEG Technical Program Expanded Abstracts, 30:2767–2771, 2011.
- Nicholas Cory Williams. Geologically-Constrained UBC-GIF Gravity and Magnetic Inversions with Examples from the Agnew-Wiluna Greenstone Belt, Western Australia. PhD thesis, University of British Columbia, 2008.