

J. Wolberg

# Data Analysis Using the Method of Least Squares

Extracting the Most  
Information from Experiments

 Springer

J. Wolberg

---

# Data Analysis Using the Method of Least Squares

Extracting the Most Information from Experiments

With 58 Figures and 68 Tables



Springer

John Wolberg

Technion-Israel Institute of Technology  
Faculty of Mechanical Engineering  
32000 Haifa, Israel  
E-mail: jwolber@attglobal.net

Library of Congress Control Number: 2005934230

ISBN-10 3-540-25674-1 Springer Berlin Heidelberg New York

ISBN-13 978-3-540-25674-8 Springer Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

Springer is a part of Springer Science+Business Media.

springer.com

© Springer-Verlag Berlin Heidelberg 2006

Printed in Germany

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Typesetting: Data prepared by the Author and by SPI Publisher Services

Cover design: *design & production* GmbH, Heidelberg

Printed on acid-free paper SPIN 11010197 62/3141/SPI Publisher Services 5 4 3 2 1 0

*For my parents, Sidney and Beatrice Wolberg ז"ל*

*My wife Laurie*

*My children and their families:*

*Beth, Gilad, Yoni and Maya Sassoon*

*David, Pazit and Sheli Wolberg*

*Danny, Iris, Noa, Adi and Liat Wolberg*

*Tamar, Ronen, Avigail and Aviv Kimchi*

## Preface

Measurements through quantitative experiments are one of the most fundamental tasks in all areas of science and technology. Astronomers analyze data from asteroid sightings to predict orbits. Computer scientists develop models for recognizing spam mail. Physicists measure properties of materials at low temperatures to understand superconductivity. Materials engineers study the reaction of materials to varying load levels to develop methods for prediction of failure. Chemical engineers consider reactions as functions of temperature and pressure. The list is endless. From the very small-scale work on DNA to the huge-scale study of black holes, quantitative experiments are performed and the data must be analyzed.

Probably the most popular method of analysis of the data associated with quantitative experiments is least squares. It has been said that the method of least squares was to statistics what calculus was to mathematics. Although the method is hardly mentioned in most engineering and science undergraduate curricula, many graduate students end up using the method to analyze the data gathered as part of their research. There is not a lot of available literature on the subject. Very few books deal with least squares at the level of detail that the subject deserves. Many books on statistics include a chapter on least squares but the treatment is usually limited to the simplest cases of linear least squares. The purpose of this book is to fill the gaps and include the type of information helpful to scientists and engineers interested in applying the method in their own special fields.

The purpose of many engineering and scientific experiments is to determine parameters based upon a mathematical model related to the phenomenon under observation. Even if the data is analyzed using least squares, the full power of the method is often overlooked. For example, the data can be weighted based upon the estimated errors associated with the data. Results from previous experiments or calculations can be combined with the least squares analysis to obtain improved estimate of the model parameters. In addition, the results can be used for predicting values of the dependent variable or variables and the associated uncertainties of the predictions as functions of the independent variables.

The introductory chapter (Chapter 1) includes a review of the basic statistical concepts that are used throughout the book. The method of least squares is developed in Chapter 2. The treatment includes development of mathematical models using both linear and nonlinear least squares. In Chapter 3 evaluation of models is considered. This chapter includes methods for measuring the "goodness of fit" of a model and methods for comparing different models. The subject of candidate predictors is discussed in Chapter 4. Often there are a number of candidate predictors and the task of the analyst is to try to extract a model using subspaces of the full candidate predictor space. In Chapter 5 attention is turned towards designing experiments that will eventually be analyzed using least squares. The subject considered in Chapter 6 is nonlinear least squares software. Kernel regression is introduced in the final chapter (Chapter 7). Kernel regression is a nonparametric modeling technique that utilizes local least squares estimates.

Although general purpose least squares software is available, the subject of least squares is simple enough so that many users of the method prefer to write their own routines. Often, the least squares analysis is a part of a larger program and it is useful to imbed it within the framework of the larger program. Throughout the book very simple examples are included so that the reader can test his or her own understanding of the subject. These examples are particularly useful for testing computer routines.

The REGRESS program has been used throughout the book as the primary least squares analysis tool. REGRESS is a general purpose nonlinear least squares program and I am its author. The program can be downloaded from [www.technion.ac.il/wolberg](http://www.technion.ac.il/wolberg).

I would like to thank David Aronson for the many discussions we have had over the years regarding the subject of data modeling. My first experiences with the development of general purpose nonlinear regression software were influenced by numerous conversations that I had with Marshall Rafal. Although a number of years have passed, I still am in contact with Marshall. Most of the examples included in the book were based upon software that I developed with Ronen Kimchi and Victor Leikehman and I would like to thank them for their advice and help. I would like to thank Ellad Tadmor for getting me involved in the research described in Section 7.7. Thanks to Richard Green for introducing me to the first English translation of Gauss's *Theoria Motus* in which Gauss developed the foundations of the method of least squares. I would also like to thank Donna Bossin for her help in editing the manuscript and teaching me some of the cryptic subtleties of WORD.

I have been teaching a graduate course on analysis and design of experiments and as a result have had many useful discussions with our students throughout the years. When I decided to write this book two years ago, I asked each student in the course to critically review a section in each chapter that had been written up to that point. Over 20 students in the spring of 2004 and over 20 students in the spring of 2005 submitted reviews that included many useful comments and ideas. A number of typos and errors were located as a result of their efforts and I really appreciated their help.

John R. Wolberg  
Haifa, Israel  
July, 2005

# Contents

<b>Chapter 1</b>	<b>INTRODUCTION .....</b>	<b>1</b>
1.1	Quantitative Experiments.....	1
1.2	Dealing with Uncertainty .....	5
1.3	Statistical Distributions .....	6
	The normal distribution.....	8
	The binomial distribution .....	10
	The Poisson distribution.....	11
	The $\chi^2$ distribution.....	13
	The $t$ distribution .....	15
	The $F$ distribution.....	16
1.4	Parametric Models .....	17
1.5	Basic Assumptions .....	19
1.6	Systematic Errors .....	22
1.7	Nonparametric Models .....	24
1.8	Statistical Learning .....	27
<b>Chapter 2</b>	<b>THE METHOD OF LEAST SQUARES .....</b>	<b>31</b>
2.1	Introduction.....	31
2.2	The Objective Function.....	34
2.3	Data Weighting .....	38



2.4	Obtaining the Least Squares Solution.....	44
2.5	Uncertainty in the Model Parameters.....	50
2.6	Uncertainty in the Model Predictions .....	54
2.7	Treatment of Prior Estimates .....	60
2.8	Applying Least Squares to Classification Problems .....	64
<b>Chapter 3 MODEL EVALUATION.....</b>		<b>73</b>
3.1	Introduction.....	73
3.2	Goodness-of-Fit .....	74
3.3	Selecting the Best Model .....	79
3.4	Variance Reduction.....	85
3.5	Linear Correlation.....	88
3.6	Outliers .....	93
3.7	Using the Model for Extrapolation .....	96
3.8	Out-of-Sample Testing .....	99
3.9	Analyzing the Residuals .....	105
<b>Chapter 4 CANDIDATE PREDICTORS.....</b>		<b>115</b>
4.1	Introduction.....	115
4.2	Using the $F$ Distribution .....	116
4.3	Nonlinear Correlation .....	122
4.4	Rank Correlation.....	131
<b>Chapter 5 DESIGNING QUANTITATIVE EXPERIMENTS.....</b>		<b>137</b>
5.1	Introduction.....	137
5.2	The Expected Value of the Sum-of-Squares.....	139
5.3	The Method of Prediction Analysis .....	140
5.4	A Simple Example: A Straight Line Experiment.....	143
5.5	Designing for Interpolation.....	147
5.6	Design Using Computer Simulations.....	150
5.7	Designs for Some Classical Experiments .....	155
5.8	Choosing the Values of the Independent Variables .....	162

5.9	Some Comments about Accuracy .....	167
<b>Chapter 6</b>	<b>SOFTWARE .....</b>	<b>169</b>
6.1	Introduction.....	169
6.2	General Purpose Nonlinear Regression Programs .....	170
6.3	The NIST Statistical Reference Datasets .....	173
6.4	Nonlinear Regression Convergence Problems.....	178
6.5	Linear Regression: a Lurking Pitfall .....	184
6.6	Multi-Dimensional Models .....	191
6.7	Software Performance.....	196
6.8	The REGRESS Program .....	198
<b>Chapter 7</b>	<b>KERNEL REGRESSION .....</b>	<b>203</b>
7.1	Introduction.....	203
7.2	Kernel Regression Order Zero .....	205
7.3	Kernel Regression Order One.....	208
7.4	Kernel Regression Order Two .....	212
7.5	Nearest Neighbor Searching .....	215
7.6	Kernel Regression Performance Studies.....	223
7.7	A Scientific Application .....	225
7.8	Applying Kernel Regression to Classification.....	232
7.9	Group Separation: An Alternative to Classification .....	236
<b>Appendix A:</b>	<b>Generating Random Noise .....</b>	<b>239</b>
<b>Appendix B:</b>	<b>Approximating the Standard Normal Distribution .....</b>	<b>243</b>
<b>References</b>	<b>.....</b>	<b>245</b>
<b>Index</b>	<b>.....</b>	<b>249</b>

# Chapter 1 INTRODUCTION

## 1.1 Quantitative Experiments

Most areas of science and engineering utilize **quantitative experiments** to determine parameters of interest. Quantitative experiments are characterized by measured variables, a mathematical model and unknown parameters. For most experiments the method of **least squares** is used to analyze the data in order to determine values for the unknown parameters.

As an example of a quantitative experiment, consider the following: measurement of the half-life of a radioactive isotope. Half-life is defined as the time required for the count rate of the isotope to decrease by one half. The experimental setup is shown in Figure 1.1.1. Measurements of **Counts** (i.e., the number of counts observed per time unit) are collected from time 0 to time **tmax**. The mathematical model for this experiment is:

$$\text{Counts} = \text{amplitude} \cdot e^{-\text{decay\_constant} \cdot t} + \text{background} \quad (1.1.1)$$

For this experiment, **Counts** is the **dependent variable** and time **t** is the **independent variable**. For this mathematical model there are 3 unknown parameters (**amplitude**, **decay\_constant** and **background**). Possible sources of the background "noise" are cosmic radiation, noise in the instrumentation and sometimes a second much longer lived radioisotope within the source. The analysis will yield values for all three parameters but only the value of **decay\_constant** is of interest. The half-life is determined from the resulting value of the decay constant:

$$e^{-\text{decay\_constant} \cdot \text{half\_life}} = 1/2$$

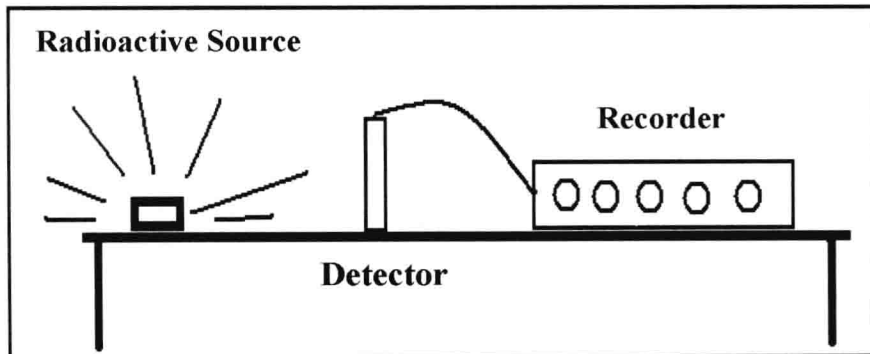
$$\text{half\_life} = \frac{0.69315}{\text{decay\_constant}} \quad (1.1.2)$$

The number 0.69315 is the natural logarithm of 2. This mathematical model is based upon the physical phenomenon being observed: the number of counts recorded per unit time from the radioactive isotope decreases exponentially to the point where all that is observable is the background noise.

There are alternative methods for conducting and analyzing this experiment. For example, the value of *background* could be measured in a separate experiment. One could then subtract this value from the observed values of *Counts* and then use a mathematical model with only two unknown parameters (*amplitude* and *decay\_constant*):

$$\text{Counts} - \text{background} = \text{amplitude} \cdot e^{-\text{decay\_constant} \cdot t} \quad (1.1.3)$$

The selection of a mathematical model for a particular experiment might be trivial or it might be the main thrust of the work. Indeed, the purpose of many experiments is to either prove or disprove a particular mathematical model. If, for example, a mathematical model is shown to agree with experimental results, it can then be used to make predictions of the dependent variable for other values of the independent variables.



**Figure 1.1.1 Experiment to Measure Half-life of a Radioisotope**

Another important aspect of experimental work relates to the determination of the unknown parameters. Besides evaluation of these parameters by experiment, there might be an alternative calculation of the parameters based upon theoretical considerations. The purpose of the experiments for such cases is to confirm the theoretical results. Indeed, experiments go hand-in-hand with theory to improve our knowledge of the world around us.

Equations (1.1.1) and (1.1.3) are examples of mathematical models with only one independent variable (i.e., time  $t$ ) and only one dependent variable (i.e., **Counts**). Often the mathematical model requires several independent variables and sometimes even several dependent variables. For example, consider classical chemical engineering experiments in which reaction rates are measured as functions of both pressure and temperature:

$$\text{reaction\_rate} = f(\text{pressure}, \text{temperature}) \quad (1.1.4)$$

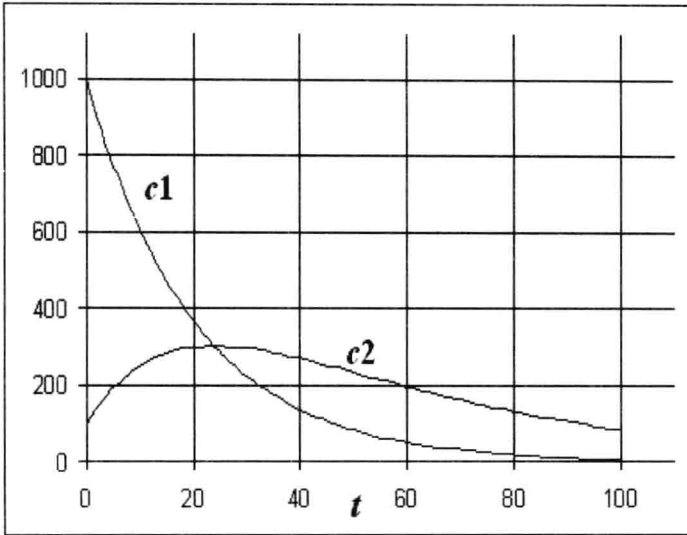
The actual form of the function  $f$  is dependent upon the type of reaction being studied.

The following example relates to an experiment that requires two dependent variables. This experiment is a variation of the experiment illustrated in Figure 1.1.1. Some radioactive isotopes decay into a second radioisotope. The decays from both isotopes give off signals of different energies and appropriate instrumentation can differentiate between the two different signals. We can thus measure count rates from each isotope simultaneously. If we call them  $c1$  and  $c2$ , assuming background radiation is negligible, the appropriate mathematical model would be:

$$c1 = a1 \cdot e^{-d1 \cdot t} \quad (1.1.5)$$

$$c2 = a2 \cdot e^{-d2 \cdot t} + a1 \frac{d2}{d2 - d1} (e^{-d1 \cdot t} - e^{-d2 \cdot t}) \quad (1.1.6)$$

This model contains four unknown parameters: the two amplitudes ( $a1$  and  $a2$ ) and the two decay constants ( $d1$  and  $d2$ ). The two dependent variables are  $c1$  and  $c2$ , and the single independent variable is time  $t$ . The time dependence of  $c1$  and  $c2$  are shown in Figure 1.1.2 for one set of the parameters.



**Figure 1.1.2** Counts versus Time for Equations 1.1.5 and 1.1.6  
 $a1=1000, a2=100, d1=0.05, d2=0.025$

The purpose of conducting experiments is not necessarily to prove or disprove a mathematical model or to determine parameters of a model. For some experiments the only purpose is to extract an equation from the data that can be used to predict values of the dependent variable (or variables) as a function of the independent variable (or variables). For such experiments the data is analyzed using different proposed equations (i.e., mathematical models) and the results are compared in order to select a "best" model.

We see that there are different reasons for performing quantitative experiments but what is common to all these experiments is the task of data analysis. In fact, there is no need to differentiate between physical experiments and experiments based upon computer generated data. Once data has been obtained, regardless of its origin, the task of data analysis commences. Whether or not the method of least squares is applicable depends upon the applicability of some basic assumptions. A discussion of the conditions allowing least squares analysis is included in Section 1.5: **Basic Assumptions.**

## 1.2 Dealing with Uncertainty

The estimation of uncertainty is an integral part of data analysis. It is not enough to just measure something. We always need an estimate of the accuracy of our measurements. For example, when we get on a scale in the morning, we know that the uncertainty is plus or minus a few hundred grams and this is considered acceptable. If, however, our scale were only accurate to plus or minus 10 kilograms this would be unacceptable. For other measurements of weight, an accuracy of a few hundred grams would be totally unacceptable. For example, if we wanted to purchase a gold bar, our accuracy requirements for the weight of the gold bar would be much more stringent. When performing quantitative experiments, we must take into consideration uncertainty in the input data. Also, the output of our analysis must include estimates of the uncertainty of the results. One of the most compelling reasons for using least squares analysis of data is that uncertainty estimates are obtained quite naturally as a part of the analysis. For almost all applications the standard deviation ( $\sigma$ ) is the accepted measure of uncertainty. Let us say we need an estimate of the uncertainty associated with the measurement of the weight of gold bars. One method for obtaining such an estimate is to repeat the measurement  $n$  times and record the weights  $w_i$ ,  $i = 1$  to  $n$ . The estimate of  $\sigma$  (the estimated standard deviation of the weight measurement) is computed as follows:

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (w_i - w_{avg})^2 \quad (1.2.1)$$

In this equation  $w_{avg}$  is the average value of the  $n$  measurements of  $w$ . The need for  $n-1$  in the denominator of this equation is best explained by considering the case in which only one measurement of  $w$  is made (i.e.,  $n = 1$ ). For this case we have no information regarding the "spread" in the measured values of  $w$ .

Fortunately, for most measurements we don't have to estimate  $\sigma$  by repeating the measurement many times. Often the instrument used to perform the measurement is provided with some estimation of the accuracy of the measurements. Typically the estimation of  $\sigma$  is provided as a fixed percentage (e.g.,  $\sigma = 1\%$ ) or a fixed value (e.g.,  $\sigma = 0.5$  grams). Sometimes the accuracy is dependent upon the value of the quantity being measured in a more complex manner than just a fixed percentage or a constant value. For such cases the provider of the measurement instrument might supply

this information in a graphical format or perhaps as an equation. For cases in which the data is calculated rather than measured, the calculation is incomplete unless it is accompanied by some estimate of uncertainty.

Once we have an estimation of  $\sigma$ , how do we interpret it? In addition to  $\sigma$ , we have a result either from measurements or from a calculation. Let us define the result as  $x$  and the true (but unknown value) of what we are trying to measure or compute as  $\mu$ . Typically we assume that our best estimate of this true value of  $\mu$  is  $x$  and that  $\mu$  is located within a region around  $x$ . The size of the region is characterized by  $\sigma$ . A typical assumption is that the probability of  $\mu$  being greater or less than  $x$  is the same. In other words, our measurement or calculation includes a random error characterized by  $\sigma$ . Unfortunately this assumption is not always valid!

Sometimes our measurements or calculations are corrupted by **systematic errors**. Systematic errors are errors that cause us to either systematically under-estimate or over-estimate our measurements or computations. One source of systematic errors is an unsuccessful calibration of a measuring instrument. Another source is failure to take into consideration external factors that might affect the measurement or calculation (e.g., temperature effects). Data analysis of quantitative experiments is based upon the assumption that the measured or calculated independent and dependent variables are not subject to systematic errors. If this assumption is not true, then errors are introduced into the results that do not show up in the computed values of the  $\sigma$ 's. One can modify the least squares analysis to study the sensitivity of the results to systematic errors but whether or not systematic errors exist is a fundamental issue in any work of an experimental nature.

### 1.3 Statistical Distributions

In nature most quantities that are observed are subject to a statistical distribution. The distribution is often inherent in the quantity being observed but might also be the result of errors introduced in the method of observation. An example of an inherent distribution can be seen in a study in which the percentage of smokers is to be determined. Let us say that one thousand people above the age of 18 are tested to see if they are smokers. The percentage is determined from the number of positive responses. It is obvious that if 1000 different people are tested the result will be different. If many groups of 1000 were tested we would be in a position to say some-



thing about the distribution of this percentage. But do we really need to test many groups? Knowledge of statistics can help us estimate the standard deviation of the distribution by just considering the first group!

As an example of a distribution caused by a measuring instrument, consider the measurement of temperature using a thermometer. Uncertainty can be introduced in several ways:

- 1) The persons observing the result of the thermometer can introduce uncertainty. If, for example, a nurse observes a temperature of a patient as  $37.4^{\circ}\text{C}$ , a second nurse might record the same measurement as  $37.5^{\circ}\text{C}$ . (Modern thermometers with digital outputs can eliminate this source of uncertainty.)
- 2) If two measurements are made but the time taken to allow the temperature to reach equilibrium is different, the results might be different. (Taking care that sufficient time is allotted for the measurement can eliminate this source of uncertainty.)
- 3) If two different thermometers are used, the instruments themselves might be the source of a difference in the results. This source of uncertainty is inherent in the quality of the thermometers. Clearly, the greater the accuracy, the higher is the quality of the instrument and usually, the greater the cost. It is far more expensive to measure a temperature to  $0.001^{\circ}\text{C}$  than  $0.1^{\circ}\text{C}$ !

We use the symbol  $\Phi$  to denote a distribution. Thus  $\Phi(x)$  is the distribution of some quantity  $x$ . If  $x$  is a discrete variable then the definition of  $\Phi(x)$  is:

$$\sum_{xmin}^{xmax} \Phi(x) = 1 \quad (1.3.1)$$

If  $x$  is a continuous variable:

$$\int_{xmin}^{xmax} \Phi(x) dx = 1 \quad (1.3.2)$$