

Determination of Scaling Laws from Statistical Data

Patricio Mendez* and Fernando Ordóñez†

Abstract

Providing a simple explanation of an engineering process typically requires extensive knowledge of the process itself. We present in this paper a statistically based methodology that obtains a scaling factor and a sequence of dimensionless groups of increasing relevance for a given process of interest. The input to our algorithm is experimental data of the process and information of the units for the variables involved. These results enable one to simplify the analysis of the process by considering only the most relevant parameters, without requiring any specific knowledge of the process in question. We apply this methodology to the problem of ceramic to metal joining for illustration, and compare the our algorithm with results available in the literature.

1 Introduction

Engineering processes involve several degrees of freedom; however, experienced people can often understand and control these processes by considering a few of these parameters. For example, when studying some types of ceramics to metals joining, elastic effects are included and plastic effects are excluded to build a simple explanation of the process [1]. The process of selection of relevant variables is currently done by making “educated guesses,” and this is the reason why experience is so important for the task.

In this paper we present a procedure to identify relevant parameters from the analysis of the experimental data of a process. This procedure combines a linear regression model that incorporates this experimental data, and physical considerations of the process, namely that the units of resulting model match the units of the dependent variable. We denote this additional constraint imposed on the linear regression model the *units constraint*. We therefore look for the model that minimizes the prediction error among the models that have the correct units. For this approach to be well defined, we require that the linear regression model be able to satisfy the units constraint.

Our approach differs from that presented by Vignaux [2, 3] in that the dimensionless groups are determined from the data instead of being created *a priori*. In

*Exponent, 21 Strathmore Road, Natick, MA 01760, USA, email: pmendez@exponent.com

†Industrial and Systems Engineering, University of Southern California, GER-247, Los Angeles, CA 90089-0193, USA, email: fordon@usc.edu

the procedure presented here, once we obtain a linear model that satisfies the units constraint, we decompose the expression into several dimensionless groups of increasing relevance. This way, we establish a bridge with dimensional analysis by providing a base of dimensionless groups for the system. Furthermore, we reach beyond the capabilities of dimensional analysis by determining the relevance of the dimensionless groups. This way the problem can be simplified by considering only the most relevant parameters, just like experienced engineers do.

The current algorithm provides a scaling factor and a set of dimensionless groups ordered by their relevance to the problem. The scaling factor can be corrected using the dimensionless groups until a desired balance between accuracy and simplicity is achieved. In effect, this algorithm provides a way to regroup the degrees of freedom of a problem in combinations that greatly decrease the number of degrees of freedom necessary to reproduce the experimental data.

The dimensionless groups have utility beyond the database used to obtain them. Neglecting the dimensionless groups of lesser relevance we reduce the number of relevant parameters in the problem; this way we can design experiments that require fewer experimental points.

2 The Engineering Approximation

We assume that the physical process being studied involves n parameters (or variables), which we label X_1, \dots, X_n . We also postulate that a property of this process that is of interest to us, Y , is determined by a physical law in the form of a power law. In other words, we assume that Y is completely described by the following equation:

$$Y = Y_S \prod_{i=1}^m \Pi_i , \quad (1)$$

where the term Y_S is a scaling factor, which has the same units as Y , and the Π_i are m dimensionless groups. We choose this power law model for three reasons: 1) The combination of units has the form of a power law. 2) The expressions of many physical phenomena have the form of power laws. 3) Many empirical regressions of engineering data in log-log plots tend to give a straight line, which corresponds to a power law.

The power law model also implies that the scaling factor and the dimensionless groups depend on the parameters (variables) in a product form, that is

$$Y_S = a \prod_{j=1}^n X_j^{a_{0j}} , \quad \Pi_i = \prod_{j=1}^n X_j^{a_{ij}} , \quad (2)$$

where the term a is a numerical constant. Incorporating the product form of the scaling factor and the dimensionless groups into equation (1), we obtain the following expression describing the process:

$$Y = a \prod_{j=1}^n X_j^{a_{0j}} \prod_{i=1}^m \prod_{j=1}^n X_j^{a_{ij}} ,$$

which rearranging becomes

$$Y = a \prod_{j=1}^n X_j^{\sum_{i=0}^m a_{ij}} . \quad (3)$$

The objective of our methodology is to obtain the expression given by (1) from the experimental data of the process. To recover expression (1), we have to identify the scaling factor Y_S and the dimensionless groups Π_i in the order of significance to the dependent variable Y . We outline the central ideas of this methodology in Section 3.

An additional consideration that the methodology will make use of, is that the exponents in equation (3) are such that the quantity on the right has the same units as the dependent variable Y . This fact remains true even after eliminating dimensionless groups. This additional constraint, which we call the units constraint, can be expressed by

$$\text{units of } Y = a \prod_{j=1}^n (\text{units of } X_j)^{\sum_{i=0}^m a_{ij}} . \quad (4)$$

The units constraint will be enforced explicitly in the methodology.

Similarly as it happens in dimensional analysis, the enforcement of the units constraint requires that our database from experiments includes measurements of all relevant parameters to build a model. Forgetting to include a parameter can have serious consequences, ranging from an incorrect model, to the impossibility of building a model at all.

In Section 4 we present the algorithm. The results this methodology obtains in a ceramic to metal joining example are presented in Section 5. In Section 6 we present a discussion of this methodology, and we conclude the paper summarizing our results in Section 7.

3 Methodology

The motivation for this methodology comes from the fact that taking the logarithm of expression (3) above yields the following equation, which is strikingly similar to a linear regression model:

$$\log Y = \beta_0 + \sum_{j=1}^n \beta_j \log X_j ,$$

where the coefficients are given by:

$$\beta_0 = \log a , \quad \beta_j = \sum_{i=0}^m a_{ij} .$$

Because in practice we consider only a finite number of variables, we will not consider variables that could have a tiny effect in determining the variable of interest Y . To account for this restriction of the model we incorporate a small normally

distributed error to this model to account for the effect of any variable beyond the n independent variables being considered. Therefore, the model which explains our dependent variable $\log Y$ is exactly a linear regression model:

$$\log Y = \beta_0 + \sum_{j=1}^n \beta_j \log X_j + \varepsilon , \quad (5)$$

where ε is a normally distributed random variable.

From experiments of the physical process in consideration we obtain p observations of the dependent property Y and of the n independent variables X_1, \dots, X_n . These observations allow us to obtain estimators to the linear model above (5) using the standard linear regression machinery. We denote the observations of the dependent variable Y by y_1, \dots, y_p and the observations for the j -th independent variable X_j by x_{1j}, \dots, x_{pj} . Grouping this data into matrix notation, we denote

$$\tilde{y} = \begin{pmatrix} \log y_1 \\ \vdots \\ \log y_p \end{pmatrix} , \quad \text{and} \quad \tilde{X} = \begin{bmatrix} 1 & \log x_{11} & \cdots & \log x_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \log x_{p1} & \cdots & \log x_{pn} \end{bmatrix} .$$

Given this notation, the estimate for the coefficients in model (5) that minimizes the residual sum of squares, in other words that solves the problem

$$\min_{\beta} \quad (\tilde{y} - \tilde{X}\beta)^t (\tilde{y} - \tilde{X}\beta) ,$$

is the solution to the normal equations

$$\tilde{X}^t \tilde{X} \beta = \tilde{X}^t \tilde{y} ,$$

see for example [4]. We denote this estimate by $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_n)$.

3.1 Constrained Linear Regression

The minimization of the residual sum of squares does not consider, however, the additional constraint that the resulting model has to have the same units as the dependent variable, Equation (4). Incorporating the units constraint will produce a model which can have a bigger residual sum of squares, but that now has the correct units. Note that the true model satisfies the units constraint.

Assume that the dependent variable and all the independent variables have units that are formed from q reference units (m, kg, s,...). We can therefore construct a q by $n+1$ matrix R in which the entry R_{ij} is the exponent of unit i in the units of variable X_j , for $j = 0, 1, \dots, n$. Note that we have included a variable X_0 , which accounts for the constant dimensionless term, this variable has $R_{i0} = 0$ for all i . Likewise we construct a q -dimensional vector b in which the entry b_i is the exponent of unit i in the dependent variable Y . Due to the power law model (3), the units constraint Equation (4) is equivalent to requiring that the coefficients satisfy $R\beta = b$.

The estimator of the coefficients in model (5) that minimizes the residual sum of squares and satisfies the units constraint is the solution to the problem

$$\begin{aligned} \min_{\beta} \quad & (\tilde{y} - \tilde{X}\beta)^t (\tilde{y} - \tilde{X}\beta) \\ \text{s.t.} \quad & R\beta = b . \end{aligned} \quad (6)$$

We denote by $\hat{\beta}^0$ the solution to problem (6).

3.2 Generation of Dimensionless Groups

We identify dimensionless groups in this model by removing one variable at a time while maintaining the units constraint. The difference in the models gives the dimensionless group discarded, since both models satisfy Equation (4). At each step, we will eliminate the variable for which the remaining model has the smallest residual sum of squares. This procedure is repeated until no more variables can be eliminated, because the units constraint cannot be satisfied. We now describe in more detail the step of eliminating one additional variable from the model.

Assume we have sequentially removed k variables from the model, that is the coefficients of the model satisfy an additional constraint, $M_k \beta = 0$ which forces exactly k of the coefficients $\beta_j = 0$. Let $\hat{\beta}^k$ be the solution to the linear regression, satisfying the units constraint and $M_k \beta = 0$.

To remove a variable X_j , such that $\hat{\beta}_j^k \neq 0$, from the model we simply fix the j -th coefficient $\beta_j = 0$. This additional constraint on the linear regression model will give us a linear regression that satisfies the units constraint, the constraint $M_k \beta = 0$, and in addition does not use variable X_j . Therefore, the reduced model is given as the solution to the optimization problem

$$\begin{aligned} z_k^*(j) = \min_{\beta} \quad & (\tilde{y} - \tilde{X}\beta)^t (\tilde{y} - \tilde{X}\beta) \\ \text{s.t.} \quad & R\beta = b \\ & M_k \beta = 0 \\ & \beta_j = 0. \end{aligned} \tag{7}$$

Among all variables that have $\hat{\beta}_j^k \neq 0$, we eliminate at iteration k the variable j which makes $z_k^*(j)$ smallest. The resulting model, with coefficients $\hat{\beta}^{k+1}$, best fits the given data in a least squares sense, satisfies the units constraint, and has $k+1$ coordinates equal to zero. This last constraint can be encoded with the matrix

$$M_{k+1} = \begin{bmatrix} M_k \\ e_j^t \end{bmatrix},$$

where e_j is the $n+1$ dimensional vector formed with 1 in the j -th coordinate and 0 in the other coordinates.

Note that the solutions to problem (7) in different iterations of this algorithm, say k and l with $k < l$, satisfy

$$R(\hat{\beta}^k - \hat{\beta}^l) = 0.$$

Therefore the difference $\delta^{kl} = \hat{\beta}^k - \hat{\beta}^l$ is a dimensionless vector. In addition, from the construction algorithm we see that for some coordinate j , $\hat{\beta}_j^k \neq 0$ and $\hat{\beta}_j^l = 0$, and therefore $\delta^{kl} \neq 0$. The algorithm removes at each iteration the variable which defines the dimensionless group $\delta^k = \hat{\beta}^k - \hat{\beta}^{k+1}$ that is less significant to the dependent variable Y .

The process of eliminating variables is valid while the system of constraints has a solution. Since each modification yields a dimensionless group, the number of iterations can be calculated using dimensional analysis. The methodology can perform at most

$$n + 1 - \text{rank}(R) \tag{8}$$

iterations. We obtain one dimensionless group more than using dimensional analysis because our analysis includes a dimensionless numerical constant in the regression. This constant has no influence from the point of view of units, but it contributes to smaller errors in the regressions.

4 Algorithm

We now present the algorithm that will identify the dimensionless groups of the linear regression model in order of significance to the dependent variable.

Step 0 Input: \widetilde{X} , \widetilde{y} , R , b . Set $k = 0$.

Step 1 Solve (6), let $\widehat{\beta}^0$ be the solution.

Step 2 Let $\min = \text{INF}$, $\text{ind} = -1$. For all j such that $\widehat{\beta}_j^k \neq 0$

- Solve (7).
- If $z_k^*(j) < \min$ then $\min = z_k^*(j)$, $\text{ind} = j$.

Step 3 If $\text{ind} = -1$ then STOP, all coordinates are zero. Else, let $\widehat{\beta}^{k+1}$ be the solution of problem $z_k^*(\text{ind})$. $\delta^k = \widehat{\beta}^k - \widehat{\beta}^{k+1}$.

Step 4 If $k < n - \text{rank}(R)$, then set $k = k + 1$ and goto Step 2. Else STOP, cannot eliminate more variables.

After the run of this algorithm, the dimensionless groups are defined by the coefficients $\delta^0, \delta^1, \dots, \delta^k$, in order of significance to the dependent variable. The scaling factor has coefficients given by $\widehat{\beta}^{k+1}$, which has the same units as the dependent variable, since it satisfies $R\beta = b$.

5 Example: Ceramic to Metal Joining

We will illustrate the implementation of the algorithm described above by applying it to a concrete case with a known solution: the joining of metals and ceramics.

Figure 1 shows the geometry of the problem, which consists of two semi-infinite cylinders one made of ceramic and the other of metal. These two cylinders are joined at their circular bases at high temperature. The temperature variation between the hot joining temperature and the cooler room temperature causes the ceramic and the metallic cylinder to decrease slightly in size. Typically the metallic cylinder will shrink more than the ceramic cylinder, causing very large stresses on and around the interface of the joint. These stresses weaken the joint; therefore, the calculation of these stresses is essential. The metric for these stresses is the “elastic strain energy” U (units= $\text{Pa}\cdot\text{m}^3$) accumulated in the ceramic. Scaling factors exist for cases in which the metallic cylinder behaves elastically [1]. Similar scaling factors for when the metallic cylinder experiences non-linear plasticity have been obtained only recently, by manual analysis of computational experiments [5]. In this paper we will show how the algorithm proposed obtains the almost the same scaling factor of [5], but automatically.

In this example we will model the materials properties of the metal as elastic-plastic and the ceramic as linear elastic. Therefore, the independent variables that

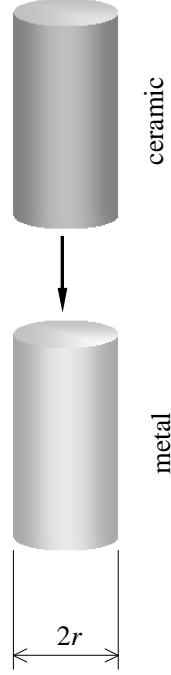


Figure 1: geometry of the ceramic and metal parts to be joined.

describe this process are the elastic modulus of the ceramic E_c (units=Pa), the elastic modulus of the metal E_m (Pa), the yield stress of the metal σ_Y (Pa), the radius r of both cylinders (m), and the differential thermal shrinkage ¹ ε_T , which is dimensionless. Table 5 lists the independent variables and the elastic strain energy, which is the dependent variable that we wish to analyze, for nine numerical simulations of ceramic to metal joints.

Since we have five independent variables ($n = 5$) and two reference units (Pa and m, which means $q = 2 = \text{rank}(R)$), our algorithm will generate four dimensionless groups ($n + 1 - \text{rank}(R)$), in addition to the scaling factor.

5.1 Output from algorithm

The successive iterations of our algorithm are shown in Table 2. This table contains the exponents of the variables for the different models and dimensionless groups. This table is separated in three sections, the top section shows the coefficients obtained in the different models, the middle section shows the dimensionless groups constructed, when applicable. The bottom section of this table reports the residual sum of squares (RSS) and the change in RSS from one model to the next.

In the top section of this table, the first column corresponds to the linear regression without the units constraint. This regression does not make physical sense because it provides a result with the wrong units. The regression indicated in the

¹ $\varepsilon_T = (\alpha_m - \alpha_c)(T_j - T_0)$, where α_c and α_m are the coefficients of thermal expansion for the ceramic and metal respectively, T_j is the joining temperature and T_0 is room temperature.

ceramic	metal	U $10^{-2}\text{Pa}\cdot\text{m}^3$	E_c 10^{11}Pa	E_m 10^{11}Pa	σ_Y 10^8Pa	r 10^{-3}m	ε_T 10^{-3}
Si_3N_4	Cu	0.423	3.04	1.28	7.58	6.25	6.85
Si_3N_4	Ni	1.52	3.04	2.08	1.48	6.25	5.15
Si_3N_4	Nb	2.80	3.04	1.03	2.40	6.25	2.10
Si_3N_4	Inco600	3.78	3.04	2.06	2.50	6.25	5.15
Si_3N_4	AISI 304	3.88	3.04	2.06	2.56	6.25	7.10
Si_3N_4	AISI 316	4.91	3.04	1.94	2.90	6.25	7.00
Al_2O_3	Ti	1.04	3.58	1.20	1.72	6.25	0.505
Al_2O_3	Inco600	3.00	3.58	2.06	2.50	6.25	2.95
Al_2O_3	AISI 304	3.16	3.58	2.00	2.56	6.25	4.90

Table 1: Input database containing the results of nine numerical experiments [5].

second column, which is only slightly less accurate, produces the correct units. This is the best regression possible that matches the units constraint.

Although the model of column two does make physical sense, it involves all the parameters of the problem, thus, it neither simplifies the problem or provides further understanding of the role of the parameters. Columns three to six show four stages of simplification. Column six provides the estimated scaling factor, which is the simplest model that fulfills the units constraint. The expression of the scaling factor is

$$\hat{U}_S = \sigma_Y r^3 \quad (9)$$

where \hat{U}_S is the estimated scaling factor.

In the middle section of the table, each of the columns three, four, five, and six yield a dimensionless group. For example, Model 6 is obtained by removing Π_4 from Model 5. Each time a dimensionless group is neglected, the new model obtained has larger residual sum of squares.

The error of the Model 6 is several orders of magnitude larger than the previous model, suggesting that it does not capture the phenomenon we are studying. Model 5 has a much smaller error, while it is only slightly more complex in its expression. Therefore, we choose Model 5 to represent our problem.

The expression of Model 5 is

$$\hat{U}_5 = \frac{\sigma_Y^{2.045} r^3}{E_c^{1.045}} \quad (10)$$

where \hat{U}_5 is the estimation of elastic strain energy using Model 5. The expression obtained manually in [5] using physical considerations is:

$$\hat{U}_{\text{ref [5]}} = \frac{\sigma_Y^2 r^3}{E_c} \quad (11)$$

Both expressions are very similar, supporting the usefulness of the algorithm proposed.

The dimensionless groups obtained with the algorithm, ranked from most to least relevant are

$$\Pi_4 = \left(\frac{\sigma_Y}{E_c} \right)^{1.045} \quad (12)$$

	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
constant	-0.424	-1.083	-0.669	0	0	0
E_c	-1.1579	-0.949	-0.816	-0.898	-1.045	0
E_m	0.194	0.173	0	0	0	0
σ_Y	1.773	1.776	1.816	1.898	2.045	1
r	2.153	3	3	3	3	3
ε_T	0.126	0.141	0.182	0.194	0	0
			Π_1	Π_2	Π_3	Π_4
constant			-0.415	-0.669	0	0
E_c			-0.133	0.082	0.147	-1.045
E_m			0.173	0	0	0
σ_Y			-0.040	-0.082	-0.147	1.045
r			0	0	0	0
ε_T			-0.041	-0.011	0.194	0
RSS	0.007	0.008	0.015	0.026	0.258	535
Δ RSS		0.001	0.007	0.011	0.232	535

Table 2: Models, dimensionless groups, and errors obtained with algorithm proposed.

$$\Pi_3 = \varepsilon_T^{0.194} \left(\frac{E_c}{\sigma_Y} \right)^{0.147} \quad (13)$$

$$\Pi_2 = e^{-0.669} \left(\frac{E_c}{\sigma_Y} \right)^{0.082} \left(\frac{1}{\varepsilon_T} \right)^{0.011} \quad (14)$$

$$\Pi_1 = e^{-0.415} \left(\frac{1}{E_c} \right)^{0.133} E_m^{0.173} \left(\frac{1}{\sigma_Y} \right)^{0.040} \left(\frac{1}{\varepsilon_T} \right)^{0.041} \quad (15)$$

6 Discussion

The most important hypothesis of our algorithm is that the process can be described by Equation (1). This is generally true when a physical phenomenon is clearly dominant over the others, for example when the metal in our example behaves almost completely as plastic. There are other power-law expressions when the metal behaves substantially as elastic. Equation (1) usually breaks down for data points that fall between the two extremes, when no phenomenon is clearly dominant. When applying the algorithm proposed, we have to be fairly sure that our data set corresponds to a situation in which the same phenomenon dominates in all of the observations.

In our example of ceramic to metal joining, we showed that the algorithm provided an expression that matched almost exactly that proposed in [5]. There are two additions to the algorithm that could make it potentially more useful to engineers, and are the subject of ongoing research.

The first of these additions would be to build an orthogonal base of dimensionless groups, that is a base of dimensionless groups in which no element of the base

contains another one in its expression. The base of dimensionless group we obtained in our example is not orthogonal, since the dimensionless group Π_3 can be expressed as $\Pi_3 = \varepsilon_T^{0.194} \Pi_4^{140}$. With this consideration, perhaps a better dimensionless group for the base would be $\Pi'_3 = \varepsilon_T^{0.194}$.

The other possible addition to the algorithm would be to define the dimensionless groups using simpler exponents, such as round numbers or rational numbers of small denominator. For example, a potentially better choice for dimensionless group would be $\Pi_4 = \sigma_Y/E_c$, this way $\Pi_4 = \Pi_4^{1.045}$.

The combination of these two additions would yield dimensionless groups of very simple expression, from which engineers can build physical interpretations. With these additions, the algorithm should provide exactly the same expression obtained in [5].

An analysis of the dimensionless groups obtained shows a difference with reference [5]. In that work, the most significant dimensionless group included the elastic modulus of the metal E_m , while the only dimensionless group that includes that parameter in our work is the least important, Π_1 . This difference deserves further study. A possible explanation of this difference is the fact that the algorithm presented here obtains one sequence of dimensionless groups and in [5] only combinations of dimensionless groups that were derived from physical reasoning were tested. Also note that the errors for removing the least important dimensionless groups (Π_1 , Π_2) are comparable. This suggests that the ranking of dimensionless groups could change with more data points.

Another topic of future research is the analysis of problems in which several variables are kept constant. This causes numerical difficulties when solving the linear regressions. Our algorithm overcomes this limitation by using the units constraint. In our example, all data points have the same radius $r = 6.25$ mm, but we could still find the proper expression because the unit constraint indicated that this parameter could only have a 3 for exponent. If we have had two constant parameters with the same units, the unit constraint would not have been able to determine which parameter to use. The same would happen if several parameters are constant but their units are not independent. This case deserves especial attention, since it can happen often in practice. For example, all the measurements might be made in the same machine which might have dimensions that cannot be changed.

7 Summary

The algorithm presented here provides a simple power-law type formula that describes a process. It also provides a set of dimensionless groups ordered by relevance. These results are of high importance to engineers that need simple formulas that bring insight during the design process. The dimensionless groups often have a physical interpretation, contributing to the understanding of the problem and the relevant factors involved.

In this work we make simultaneous use of linear regressions and dimensional analysis, building on the strengths of each other to obtain a results that could not be obtained with either of these techniques alone. Regressions alone provide expressions that have small mathematical errors, yet they are not necessary amenable to physical

interpretation. Dimensional analysis provides results physically meaningful, but cannot determine the relative importance of the dimensionless groups, this is usually done manually by experts in a given field.

The algorithm proposed has the mathematical accuracy of linear regressions and the physical meaning of dimensional analysis. It can also automatically sort the dimensionless groups, as expert engineers do, and it can overcome the problem of constant parameters, which cause singularity in regressions.

We applied this algorithm to the problem of ceramic to metal joining to obtain a simplified model of the problem. This model matched the one obtained manually in [5], showing that the algorithm can reproduce the manual process in an automatic way.

References

- [1] Bryan E. Blackwell. *A framework for determining the mechanical properties of dissimilar material joints*. Doctor of philosophy, Massachusetts Institute of Technology, Cambridge, Mass., 1996.
- [2] G. A. Vignaux and J. L. Scott. Simplifying regression models using dimensional analysis. *Australian & New Zealand Journal of Statistics*, 41(1):31–41, 1999.
- [3] G. A. Vignaux. Some examples of dimensional analysis in operations research and statistics. In *4th International Workshop on Similarity Methods*, pages 247–265, Stuttgart, Germany, 2001. University of Stuttgart.
- [4] David Freedman, Robert Pisani, and Roger Purves. *Statistics*. W.W. Norton, New York, 3rd edition, 1998.
- [5] J.-W. Park, P. F. Mendez, and T. W. Eagar. Strain energy distribution in ceramic to metal joints. *Acta Materialia*, 50:883–899, 2002.