

Support vector regression for fitting multi-variable material models

Biswajit Banerjee

Parresia Research Limited, Auckland, New Zealand

b.banerjee.nz@gmail.com

David M. Fox

Armed Forces Research Laboratory, Aberdeen Proving Ground, MD, USA

david.m.fox1.civ@mail.mil

Richard A. Regueiro

University of Colorado, Boulder, CO, USA

richard.regueiro@colorado.edu

Friday, 11 September, 2020

Abstract

Analytical expressions for phenomenological material models that depend on multiple independent variables are notoriously difficult to design. Parameter determination is also intimately tied with the model design process. Soils that exhibit elastic-plastic coupling are particularly prone to the design problem. It is not uncommon to have to redesign models for every new soil that is characterized experimentally. An unstated assumption in soil mechanics is that small inaccuracies in material models do not affect the predictive capabilities of those models significantly. First, we demonstrate that such an assumption is not warranted, particularly in the large deformation, non-monotonic loading, regime. We then proceed to explore support vector regression to replace analytical models. We show that even though support vector machines can fit input data sets accurately, they fail to generalize if the input data are overfit. Also, the approach used for scaling the input data can have a considerable effect on the quality of fits. For the small sets of data that are typically available for soils, contrary to the extant literature, we suggest a poorer fit to the input data that leads to better generalization is more robust. This approach also has the benefit that fewer support vectors are needed to model the training data. This study further emphasises the pressing need for physics-based models than can supplement phenomenological aspects of material modeling for simulations.

1 Introduction

Material test data for soils and rocks often exhibit elastic behavior that depends on the state of plastic deformation and associated internal variables. To model and simulate such materials accurately, constitutive models require that the elastic response be dependent on the plasticity state during plastic loading. For example, in the KAYENTA model (Brannon et al., 2015) experimental data are used

to fit the bulk modulus (K) model:

$$K(I_1, \varepsilon_v^p) = f_K \left[b_0 + b_1 \exp\left(-\frac{b_2}{|I_1|}\right) - b_3 \exp\left(-\frac{b_4}{|\varepsilon_v^p|}\right) \right] \quad (1)$$

where I_1 is the trace of the Cauchy stress tensor, ε_v^p is the plastic volumetric strain, f_K is a joint degradation factor, and $(b_0, b_1, b_2, b_3, b_4)$ are fitted parameters. If the material is fully or partially saturated with a fluid, these models also need to incorporate information about the porosity (ϕ), saturation level (S_w), and pore pressure (p^w). In the simplest version of the ARENA partially saturated soil model (Banerjee and Brannon, 2017; Banerjee and Brannon, 2019),

$$K(I_1^{\text{eff}}, \varepsilon_v^p, \phi, p^w) = K_d(I_1, \varepsilon_v^p) + \frac{K_f(p^w)}{\phi} \quad (2)$$

where K_d is the bulk modulus of the dry material and K_f is the bulk modulus of the fluid.

The process of determining the parameters needed for these models from multi-dimensional experimental data is nontrivial even when only two independent variables are involved. Numerous algebraic expressions have been developed to fit data because a single expression is typically unable to describe all the variations observed in experimental data.

However, nonlinear and coupled elastic-plastic models are not always thought to be necessary for accurate prediction even though some papers do discuss the issue (Homel, Guilkey, and Brannon, 2016). It is quite common for a constant bulk modulus to be fit to experimental hydrostatic loading/unloading data using optimization tools such as Dakota (Adams et al., 2009). The effect on predictions of such a choice is rarely discussed. On the other hand, nonlinear models are proposed without any discussion of the benefit of such a choice other than a better fit to the available experimental data - typically from a limited number of samples.

In this work, we first motivate the need for nonlinear bulk modulus models using a simple soil penetration simulation. In the rest of the paper we focus on support vector regression as a means of generating models that can replace algebraic expressions. e.g., equation (1). Experimental data on the hydrostatic compression loading/unloading of dry, poorly graded, concrete sand is used to illustrate the process.

Support vector machines (Cortes and Vapnik, 1995; Vapnik, 2013) have been used in numerous studies on geomaterials, but primarily for classification (Zhao and Yin, 2009; Yuvaraj et al., 2013). Regression studies are fewer in number (Xue et al., 2016; Kohestani and Hassanlourad, 2016; Zhang et al., 2017) and tend to be used to fit a small set of input data without much attention to generalizability. In this work, we explore various way in which small input data sets can be modeled with support vector regression and suggest that human judgement may be required to select a model that has adequate prediction power outside the training data set.

1.1 Variable moduli and punch impact simulation

To motivate the need for nonlinear material models, consider the plane strain approximation of a rectangular punch impacting a soil sample at high velocity. Two different soil models are simulated; one with a constant bulk modulus and the other with a pressure-dependent bulk modulus described

by equation 1 with an addition linear pressure-dependent term needed to better fit soil data. The soil is nominally equivalent to dry Colorado Mason Sand (Banerjee and Brannon, 2019).

These models have been implemented in the Material Point Method code Vaango (Banerjee and Brannon, 2017) using Arenisca (Homel, Guilkey, and Brannon, 2015). The difference between the two models can be illustrated by subjecting a single particle to a set of prescribed deformation gradients, a series of hydrostatic compression load/unload steps followed by uniaxial stress compression and release and then uniaxial tension and release. The resulting stress paths in pq -space¹ are shown in Figure 1(a). There are small differences in the stress paths and the final stress state. The differences in the two stress paths are accentuated in Figure 1(b). Though the variable modulus model achieves a higher peak stress in hydrostatic compression, the unload and uniaxial loading and unloading paths are similar. That indicates that the models are essentially identical at small strains but start to diverge as the strains increase.

Since the constant and variable modulus models produce similar results in a single particle test, we would intuit that differences would be small for more complicated situations. If we use these models to simulate the impact of a punch on soil, we observe the behavior presented in Figure 2. The punch is made of a hypoelastic steel-like material that is nominally rigid compared to the soil. The initial velocity is 200 m/s and symmetry conditions are applied to the domain boundaries, emulating plane strain conditions. At 0.9 milliseconds after impact, the stress and deformation in both the variable bulk modulus sand (left half) and the constant bulk modulus material (right half) are approximately the same. The variable modulus material has bands of higher stress and a less uniform distribution of stress. However, after the punch rebounds at the end of the impact event (at 7 milliseconds in the figure), the depth of penetration and the shear band shapes are substantially different for the two cases. This indicates that the extra effort needed to develop a nonlinear bulk modulus model may be justified for better accuracy.

2 Support vector regression

The support vector regression (SVR) approach (Schölkopf et al., 2000; Smola and Schölkopf, 2004) can be used to fit models to data without the need for closed form expressions. The advantage of this approach is that the resulting model requires few function evaluations and can, in principle, be computed as fast as a closed-form model.

For the purpose of fitting a model to the yield function, the elasticity model, or the crush curve, we can assume that the input (training) data are of the form $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\} \subset \mathbb{R}^d \times \mathbb{R}$. The aim of SVR is to find a function $y = f(\mathbf{x})$ that fits the data such that the function is as flat as possible (in $d + 1$ -dimensional space), and deviates from y_i by at most ϵ (a small quantity).

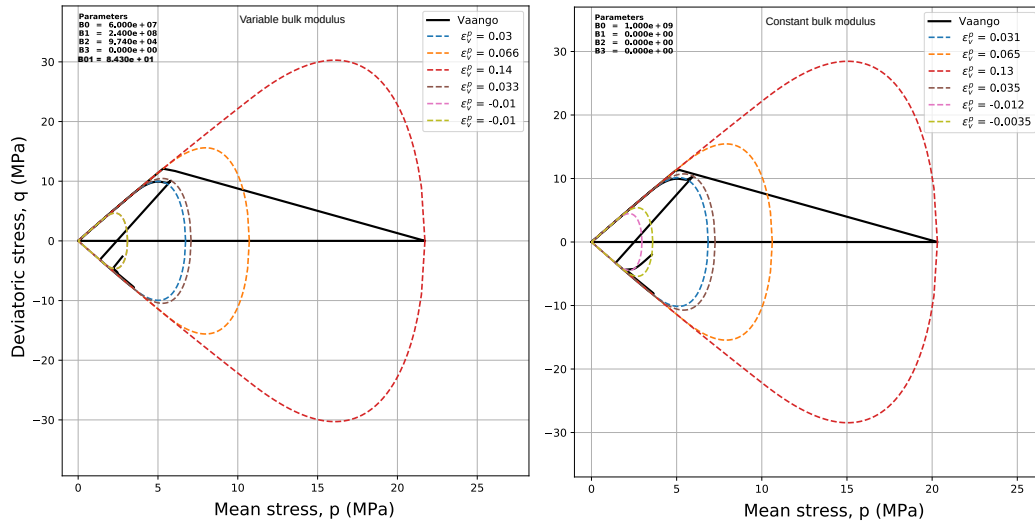
In nonlinear support vector regression we fit functions of the form

$$y = f(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x}) + b \quad (3)$$

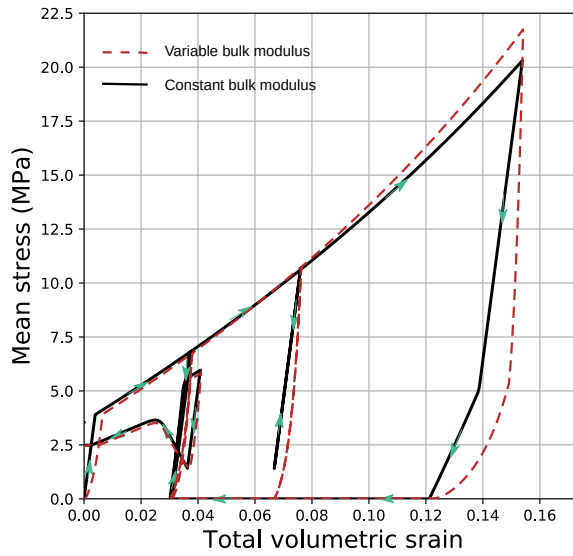
where \mathbf{w} is a vector of parameters, $\phi(\mathbf{x})$ are vector-valued basis functions, (\cdot) is an inner product, and b is a scalar offset. The fitting process can be posed as the following primal convex optimization

¹The mean stress is $p = -\frac{1}{3}\text{tr}\boldsymbol{\sigma}$ while the equivalent deviatoric stress is $q = \sqrt{3\mathbf{s} : \mathbf{s}}$ where $\mathbf{s} = \boldsymbol{\sigma} + p\mathbf{I}$ and $\boldsymbol{\sigma}$ is the Cauchy stress.

DRAFT



(a) Stress paths in pq -space. Left: Variable bulk modulus. Right: Constant bulk modulus



(b) Hydrostatic loading/unloading paths for the two models.

Figure 1 – Stress paths for single particle tests of the two material models using Vaango and Arenisca3.

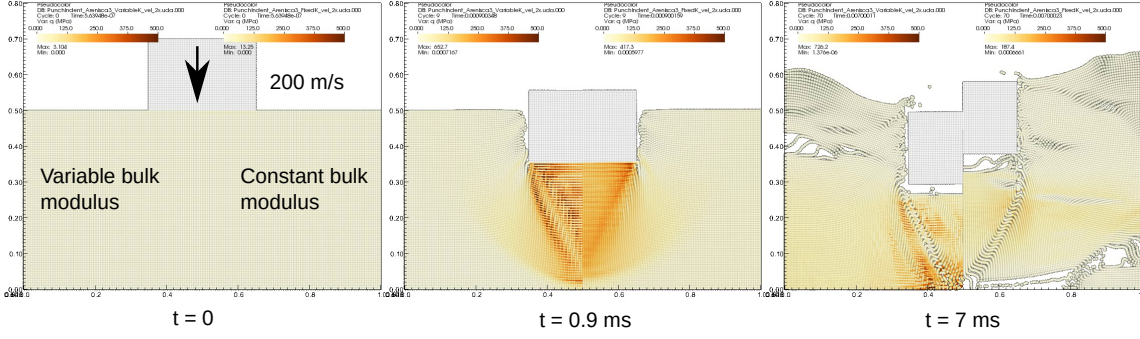


Figure 2 – Punch impacting sand at 200 m/s. Sand particles have been colored based on values of q .

problem (Vapnik, 1998):

$$\begin{aligned}
 & \underset{\mathbf{w}, b, \xi, \xi^*}{\text{minimize}} && \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\
 & \text{subject to} && \begin{cases} -(\xi_i + \epsilon) \leq y_i - \mathbf{w} \cdot \boldsymbol{\phi}(\mathbf{x}_i) - b \leq \xi_i^* + \epsilon \\ \xi_i, \xi_i^* \geq 0, \quad i = 1 \dots m \end{cases}
 \end{aligned} \tag{4}$$

where C is a constraint multiplier, m is the number of data points, and ξ_i, ξ_i^* are constraints.

In practice, it is easier to solve the dual problem for which the expansion for $f(\mathbf{x})$ becomes

$$y = f(\mathbf{x}) = \sum_{i=1}^m (\lambda_i^* - \lambda_i) K(\mathbf{x}_i, \mathbf{x}) + b, \quad K(\mathbf{x}_i, \mathbf{x}) = \boldsymbol{\phi}(\mathbf{x}_i) \cdot \boldsymbol{\phi}(\mathbf{x}) \tag{5}$$

where \mathbf{x}_i are the sample vectors, λ_i and λ_i^* are dual coefficients, and $K(\mathbf{x}_i, \mathbf{x})$ is a kernel function. The dual convex optimization problem has the form

$$\begin{aligned}
 & \underset{\lambda, \lambda^*}{\text{minimize}} && \frac{1}{2} \sum_{i,j=1}^m (\lambda_i - \lambda_i^*) K(\mathbf{x}_i, \mathbf{x}_j) (\lambda_j - \lambda_j^*) + \epsilon \sum_{i=1}^m (\lambda_i + \lambda_i^*) + \sum_{i=1}^m y_i (\lambda_i - \lambda_i^*) \\
 & \text{subject to} && \begin{cases} \sum_{i=1}^m (\lambda_i - \lambda_i^*) = 0 \\ \lambda_i, \lambda_i^* \in [0, C], \quad i = 1 \dots m. \end{cases}
 \end{aligned} \tag{6}$$

The free parameters for the fitting process are the quantities ϵ and C . SVR accuracy also depends strongly on the choice of kernel function. In this paper, we use the Gaussian radial basis function:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp \left[-\frac{(\mathbf{x}_i - \mathbf{x}_j) \cdot (\mathbf{x}_i - \mathbf{x}_j)}{d \sigma^2} \right] \tag{7}$$

where d is the dimension of \mathbf{x} and σ^2 is the width of the support of the kernel (assumed to be equal to the norm of the covariance matrix of the training data in this paper).

The minimization problem solves for the difference in the dual coefficients ($\lambda - \lambda^*$) and the intercept (b), and outputs a reduced set ($m_{SV} < m$) of values of \mathbf{x}_i called “support vectors”. Given these quantities, the function (5) can be evaluated quite efficiently, particularly if the number of support vectors

is small. SVR fits to data can be computed using software such as the LIBSVM library (Chang and Lin, 2011). A variation of the above approach, called ν -SVR (Schölkopf et al., 2000) can also be used if sufficient computational resources are available.

3 Experimental data

In this report we will attempt to fit SVR models for the bulk modulus and crush curve of a dry, poorly-graded, concrete sand described by Fox et al. (2014) and tested at the University of Maryland.² The hydrostatic loading-unloading data for that sand are shown in Figure 3(a). The loading curve, shown in green, is used to fit the crush curve model. The unloading curves are used to fit a bulk modulus model that depends on the plastic strain. Zoomed plots of these unloading curves, and the associated volumetric plastic strains, are shown in Figure 3(b).

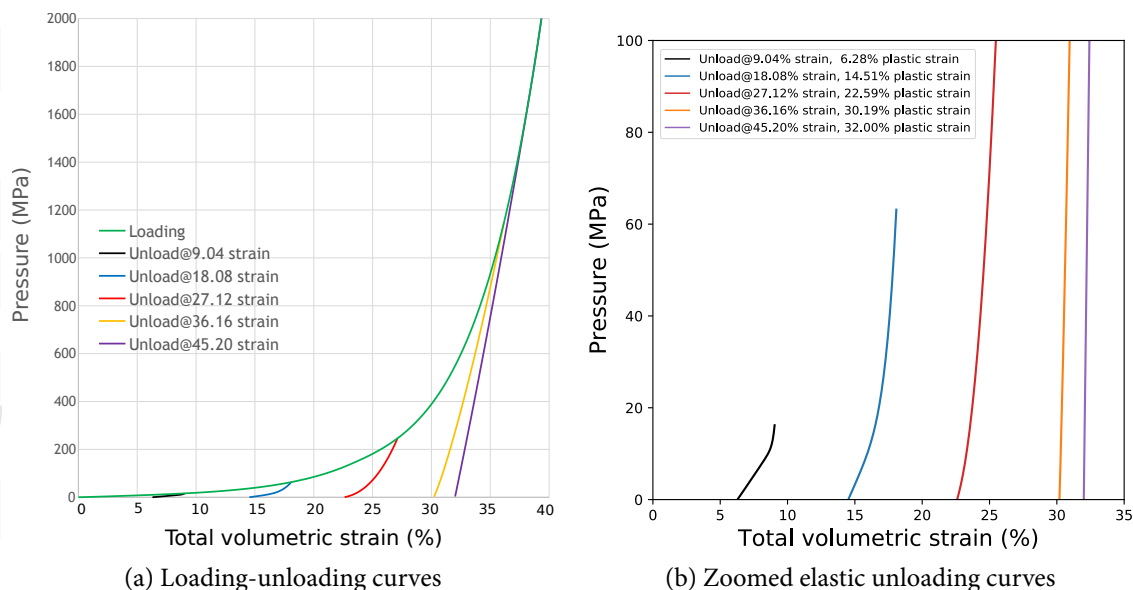


Figure 3 – Hydrostatic loading-unloading data for poorly-graded dry sand.

Strains are assumed to be additively decomposable into elastic and plastic parts. The plastic volumetric strain for each elastic unloading curve is therefore assumed to remain constant. The value of the plastic volumetric strain at for an unloading curve is determined by computing the intersection of the curve with the total strain axis. Elastic strains are computed by subtracting the plastic volumetric strain from the total strain. Figure 4(a) shows the unloading curves for the dry sand as a function of the elastic volumetric strain. Tangents to the curves represent the bulk modulus and have been plotted in Figure 4(b). to accentuate the behavior at low pressures.

Similarly, the crush-curve can be extracted from the hydrostatic compression data in Figure 3(a). The initial bulk modulus is assumed be constant (425 MPa) for states with volumetric plastic strains less than 6.28%. To compute the crush-curve pressure as a function of the plastic strain, values from

²Stephen Akers, 2018, Private communication, CCDC Army Research Laboratory, Aberdeen Proving Ground, MD, USA

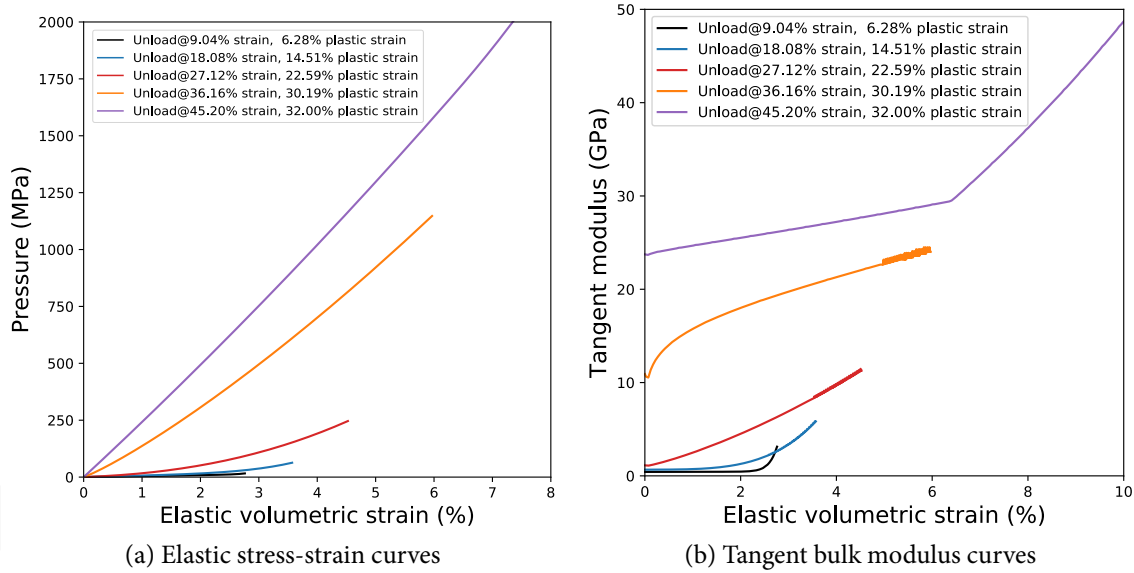


Figure 4 – Elastic unloading curves and tangent bulk modulus for poorly-graded dry sand.

the hydrostat are interpolated linearly between adjacent unloading points. Figure 5(a) depicts the hydrostat with unloading points marked with filled circles. The processed crush-curve is shown in Figure 5(b).

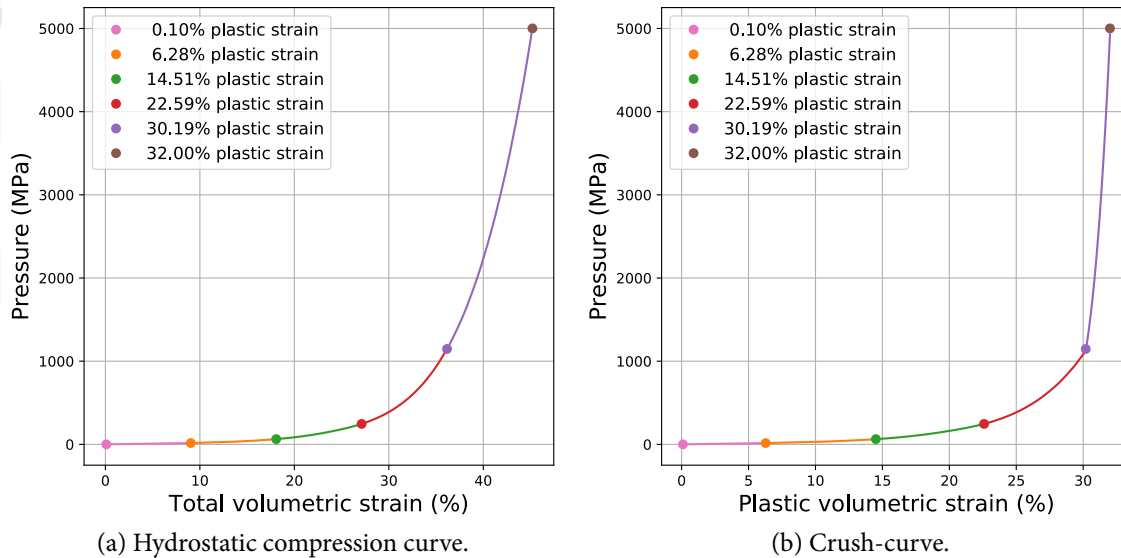


Figure 5 – Hydrostatic compression curve and crush-curve for poorly-graded dry sand.

Note that the term “crush-curve” is used more commonly to refer to the representation of the same data in a porosity versus pressure form. The porosity (ϕ) is computed using $\phi = p_3 - \varepsilon_p^v$ where ε_p^v

is the volumetric plastic strain and p_3 is the maximum volumetric plastic strain, at which all pores have been crushed (Brannon et al., 2015).

The plastic-strain dependence of the tangent bulk modulus of the sand in Figure 4 does not satisfy any obvious analytical form. It is more convenient to use a support vector regression approach to fit models to data of this nature. Also, the shapes of the stress-strain curves are more regular than the derived tangent bulk modulus curves and easier to approximate. Since we do not have unloading data for zero volumetric plastic strain, we assume that the associated stress-strain curve is identical to that for a plastic strain of 6.28%. The same assumption is made for tensile volumetric strain states, if any are present in a simulation.

4 Fitting a crush-curve model

The crush-curve depends only on the volumetric plastic strain and is therefore more straightforward to fit than the elastic unloading curves. For the regression process, the data in Figure 5(b) was read into a pandas data frame (McKinney, 2011) and scaled using the Yeo-Johnson power transform (Yeo and Johnson, 2000).³ SVR fits to the data were computed using the SVR front-end for the LIBSVM library (Chang and Lin, 2011) provided by scikit-learn (Pedregosa et al., 2011). A radial basis function kernel was used for the fitting process.

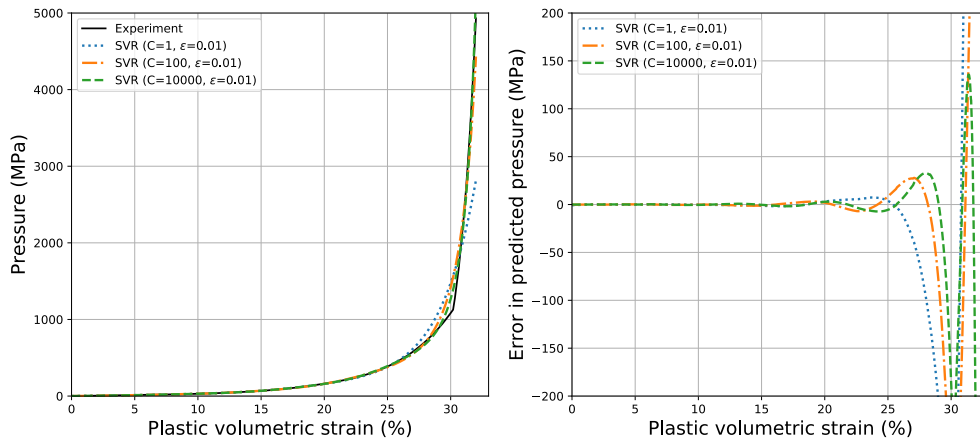
Support vector regression predictions are compared with the experimental data for fixed ϵ and varying C in Figure 6. The experimental curve is depicted with a solid line, while SVR fits are shown as dashed lines of varying color. Figure 6(a) shows fits to the pressure vs. plastic strain curve and the error in the computed pressure ($p_{\text{SVR}} - p_{\text{expt.}}$). The fits are relatively accurate up to plastic volumetric strains of around 20%. Errors increase in magnitude as the pressure increases. But, as the value of C increases, these differences decrease as a fraction of the pressure. The number of support vectors is important for the fast evaluation of the SVR model in a plasticity simulation, the smaller the better. For the data under consideration, 236 support vectors were needed for $C = 1$, while 133 were needed for $C = 100$, and 83 for $C = 10,000$.

Alternatively, if we fit the data in porosity vs. pressure form, the SVR fits and experimental curve are as shown in Figure 6(b). Note that the error is reported as a percentage of the experimental value. Visual examination of the fits immediately show that these are superior to the pressure vs. plastic strain fits. In fact, any of the three fits could be used as an approximation of the crush-curve. An artifact from the input data that shows as a discontinuous slope in the crush-curve is also handled gracefully by the SVR fits. The number of support vectors are 66, 34, and 24 for $C = 1$, 100, and 10,000, respectively. The quality of fit and the small number of support vectors suggest that using the pressure as the independent variable may be preferable when fitting crush-curves.

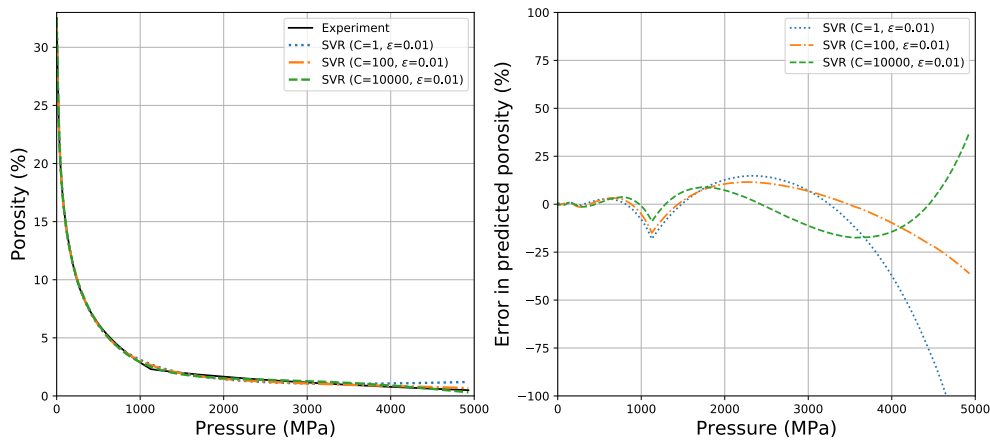
If C is kept constant at 10,000 and ϵ varied, SVR fits to the porosity-pressure crush-curve suggest that lower values of ϵ lead to better models of the experimental data. The SVR fits to the data and the percent error are shown in Figure 7. The overall quality of the fit is excellent. However, the number of support vectors are (24, 505, 570) for $\epsilon = (0.01, 0.001, 0.0001)$, respectively. Given the quality of the model for the smallest value of ϵ and also the small number of support vectors needed to achieve that accuracy, that model should be chosen for numerical simulations.

³Other transformations are possible but have been found to lead to poorer fits to the data.

DRAFT



(a) Fit to pressure-plastic strain curve.



(b) Fit to porosity-pressure curve.

Figure 6 – Support vector regression fits, at various values of C , to the crush-curve for ARL dry sand.

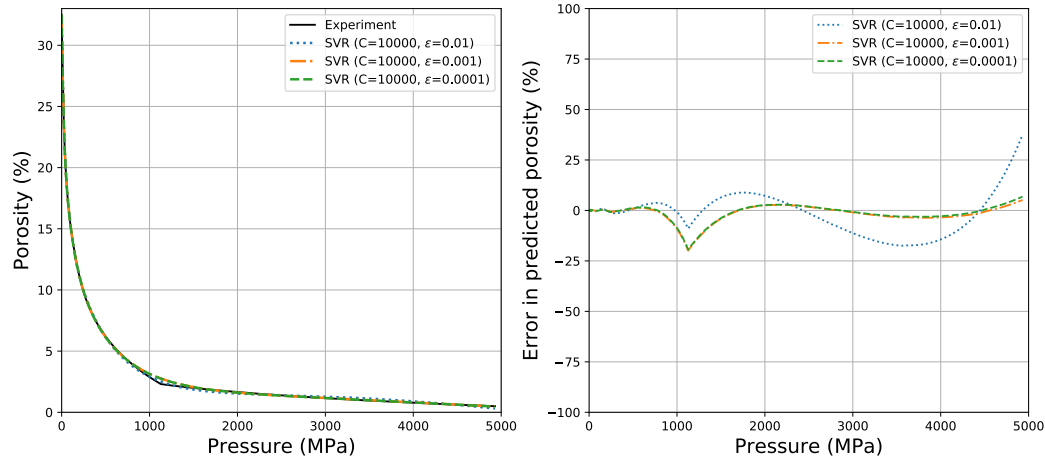


Figure 7 – Support vector regression fits, at various values of ϵ , to the porosity-pressure crush-curve for poorly-graded dry sand.

To summarize the results from this section, accurate and efficient support vector models for crush-curves can be fit if the data are expressed in porosity vs. pressure form, scaled using a power transformation, and parameters selected to optimize for both prediction accuracy and the number of support vectors.

5 Fitting a bulk modulus model

The bulk modulus for the sand in Section 3 depends on both the elastic and the plastic volumetric strain. In plasticity models, the stress and the plastic strain are typically computed before the elastic strain is known. Hence, a bulk modulus model that depends on the pressure and the plastic strain is more convenient for computations. In this paper, for convenience, we fit models for the pressure (p) as a function of elastic (ϵ_v^e) and plastic (ϵ_v^p) strains. A pressure-dependent bulk modulus model can be extracted from the $p = p(\epsilon_v^e, \epsilon_v^p)$ model by switching the independent and dependent variables.

Assuming that each unloading curve has been sampled uniformly, from Figure 3(b) it can be seen that fewer experimental data points are available at smaller plastic strains. When SVRs are fitted to these data, the curves with larger numbers of samples tend to bias the fitting process. To overcome this problem bootstrapped samples are generated where needed. Bootstrapping allows better fits to be produced in regions where physical constraints apply, such as zero pressure at zero volumetric elastic strain.

The bootstrap procedure involves resampling from existing data with or without replacement (Mooney, Duval, and Duvall, 1993) and can be used to produce samples from the underlying distribution. In the sand data set, there is only one set of unloading curves. Therefore, we repeat values at the smaller strains to increase the number of samples available for fitting. Note that increasing the number of samples can lead to a significant increase in training time and should be avoided in support vector regression if possible.

5.1 Fitting to pressure vs. elastic strain

The hydrostatic unloading data can be expressed in the form of pressure as a function of the elastic volumetric strain. Figure 4(a) shows the pressure evolution for several volumetric plastic strain. During the regression process, data for all the plastic strain values were concatenated into a single data frame and a SVR fit to the data was computed using the SVR front-end for LIBSVM from scikit-learn. Both the strains and the pressures were scaled to lie between 0 and 1 before the optimization algorithm was invoked.

If the input data are not shuffled and bootstrapping is not used, we get the fits shown in Figure 8(a) for $C = 1000$ and $\epsilon = 0.01$. The default radial basis function kernel parameter $\gamma = 1/(2\sigma^2)$ was with $\sigma^2 = \text{variance}([\epsilon_v^e, \epsilon_v^p])$. The effect on the fit of changing C was found to be minimal but decreasing the value of ϵ was found to improve the quality of the fit.⁴ From the figure (part (b)) we also see that the error in the predicted pressure increases as the elastic volumetric strain decreases.

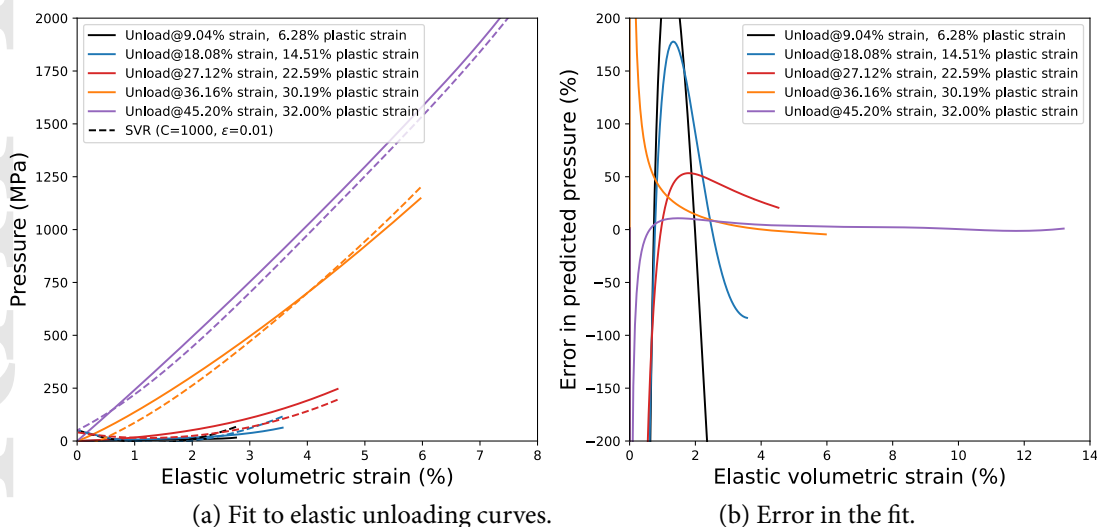


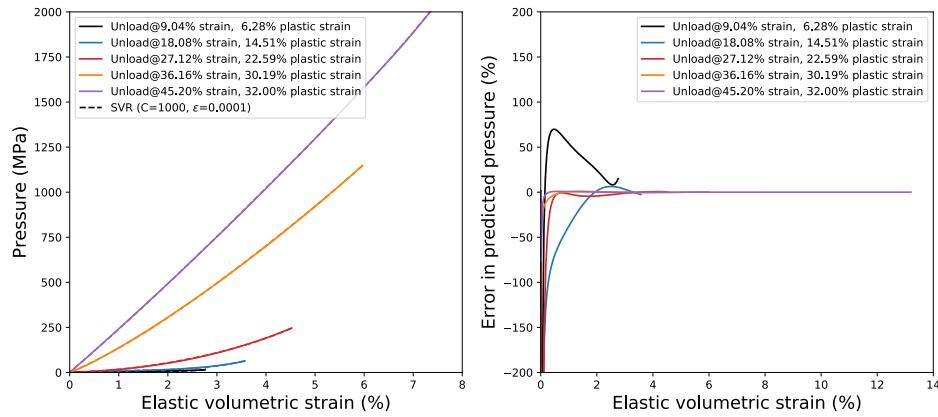
Figure 8 – Support vector regression fits ($C = 1000$, $\epsilon = 0.01$) to the original elastic unloading curves for poorly-graded dry sand.

A better fit to the input data is obtained for $\epsilon = 0.0001$ as can be seen in Figure 9(a). The effect of bootstrapping is shown in Figure 9(b), where the SVR fit used $C = 1000$ and $\epsilon = 0.0001$ and the data at the three lowest plastic strains were repeated and then the full data set was shuffled randomly. If the data are bootstrapped, the training error is typically lower for the bootstrapped samples. Also, even though large errors persist at low strains in percent terms, the absolute errors are relatively small (of the order of 10-100 MPa).

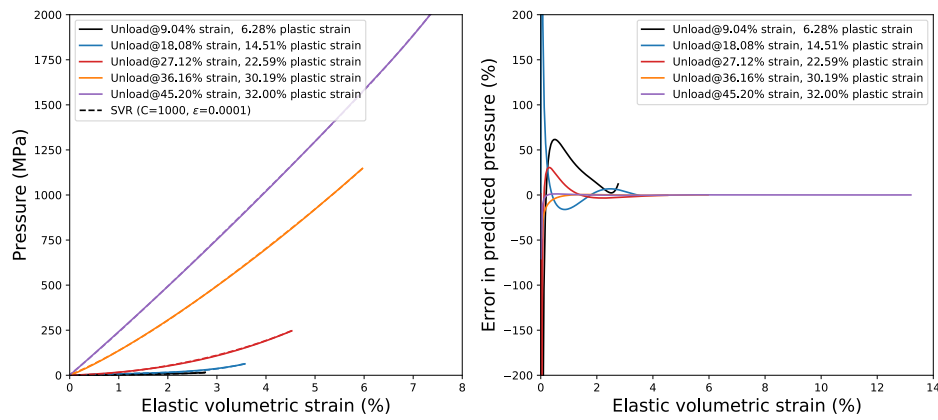
These fitted curves can be processed to extract the fitted bulk moduli. For the bootstrapped input, the fitted bulk moduli curves and the error in the fit are shown in Figure 10. The SVR fits are shown as dashed lines while the solid lines represent the experimental data. The fits increase in relative accuracy as the volumetric plastic strain increases. Also, as indicated for the previous figure, the error percentages are largest at small strains even though their magnitude is relatively small.

⁴If the pressure is not scaled, the effect of C on the quality of the fit is substantial while that of ϵ is negligible.

DRAFT



(a) Original elastic unloading curves.



(b) Bootstrapped elastic unloading curves.

Figure 9 – Support vector regression fits ($C = 1000$, $\epsilon = 0.0001$) to the original and the bootstrapped elastic unloading curves for poorly-graded dry sand.

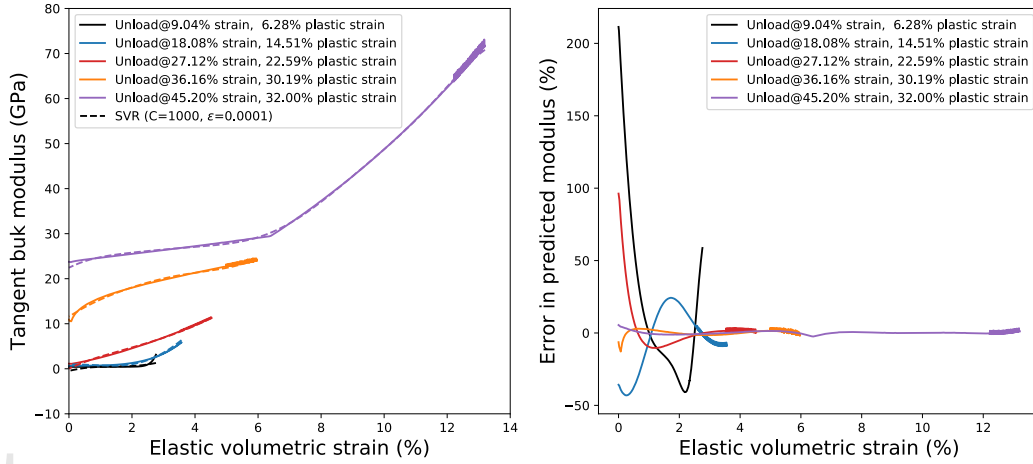


Figure 10 – Tangent bulk modulus computed from support vector regression fits ($C = 1000$, $\epsilon = 0.0001$) to bootstrapped elastic unloading curves for poorly-graded dry sand.

How well the fitted model generalize to plastic strains outside those used to train the model? Figure 11(a) shows the predictions of the SVR model trained on the bootstrapped data ($C = 1000$, $\epsilon = 0.0001$) for a range of volumetric plastic strains. The corresponding bulk modulus predictions are shown in Figure 11(b). The solid lines in the plots show the input data while the dashed lines show predicted values. From the pressure-elastic strain plots we can see that, at zero plastic strain, the predicted pressure is approximately 500 MPa while negative pressures and bulk moduli are predicted by the SVR model at 10% plastic strain. The fits improve as the plastic strain increases and the expected monotonic increase in pressure with strain holds approximately. Also, the tangent modulus curves exhibit inflection points. Clearly, even though the training data are fitted accurately by the SVR model, we cannot use the model for simulations because of its failure to generalize.

5.2 Fitting to pressure vs. total volumetric strain

The results from the previous section indicate that additional constraints, e.g., the elastic strain must be zero at zero pressure, and the pressure-strain curve at 0 plastic strain must be close to that at 6% plastic strain, are necessary for better fits to the data at small strains and better generalization. However, such constraints require the support vector regression optimizer to be modified. Since that option is not available for most users of LIBSVM and other open-source libraries, it is more convenient to examine alternative approaches to achieve the desired outcome.

One alternative is to fit support vectors to the pressure data as a function of the total volumetric strain. As before, the data are bootstrapped at plastic strains where the number of data points is small and then the data are shuffled before the fitting algorithm is invoked. Figure 12(a) shows fits to the pressure-total strain data for $C = 1000$ and $\epsilon = 0.0001$. As before, the effect of varying C is small but the fit deteriorates when ϵ is increased. The ν -SVR algorithm (Schölkopf et al., 2000) produces similar results for $\nu = 0.5$ and the same C . Cross-validation can be used to select optimal values of the control parameters. The corresponding fits to the bulk modulus are reasonable, as can be observed

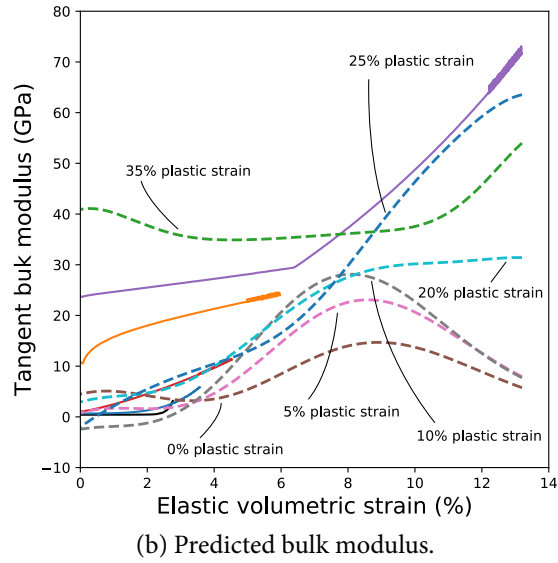
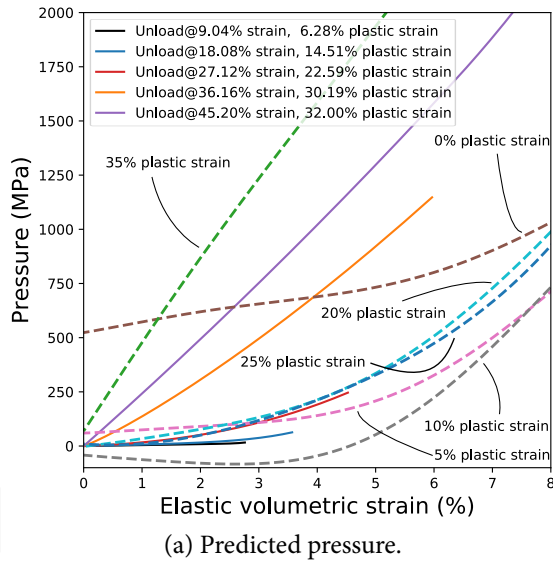


Figure 11 – Support vector regression fits ($C = 1000$, $\epsilon = 0.0001$) to the original and the bootstrapped elastic unloading curves for poorly-graded dry sand.

from Figure 12(b). However, the bulk modulus goes to zero for both 6.28% and 14.51% plastic strain due to inflections in the pressure-strain curves. Therefore, this model cannot be used in a simulation without modification.

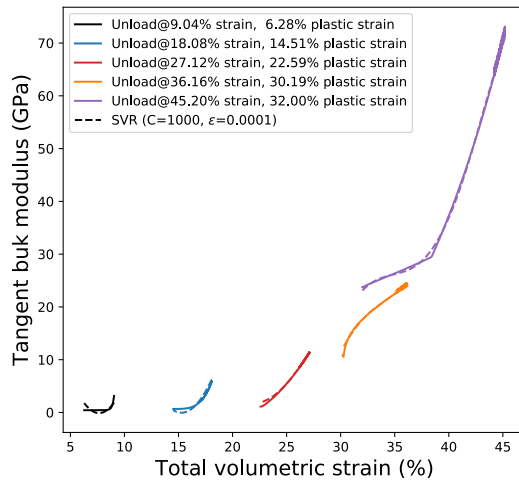
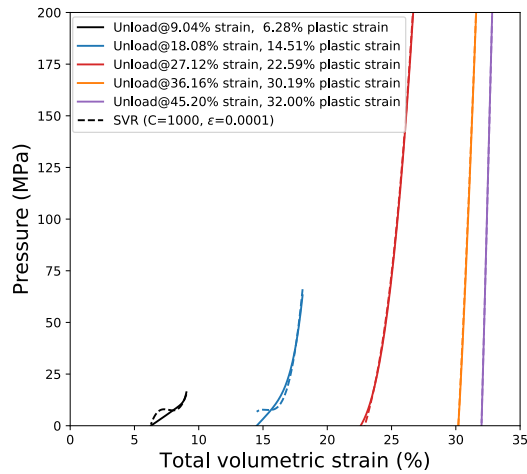


Figure 12 – Support vector regression fits ($C = 1000$, $\epsilon = 0.0001$) to the bootstrapped unloading curves for poorly-graded dry sand.

We can examine the generalizability of the fit from the predicted curves shown in Figures 13(a) and (b). The predicted purely elastic response at 0% plastic strain indicates a stress of 1250 MPa at 0 volumetric strain. The pressure decreases with increasing strain before catching up with the predicted

curve for 5% plastic strain which is also slightly negative at 5% total strain. Similar unreasonable behaviors are predicted at all the other values of plastic strain chosen for this experiment. As expected, the bulk modulus starts at a negative value or turns negative for at least two of the selected values of plastic volumetric strain. Therefore, the SVR fit not only is not very accurate but also not general.

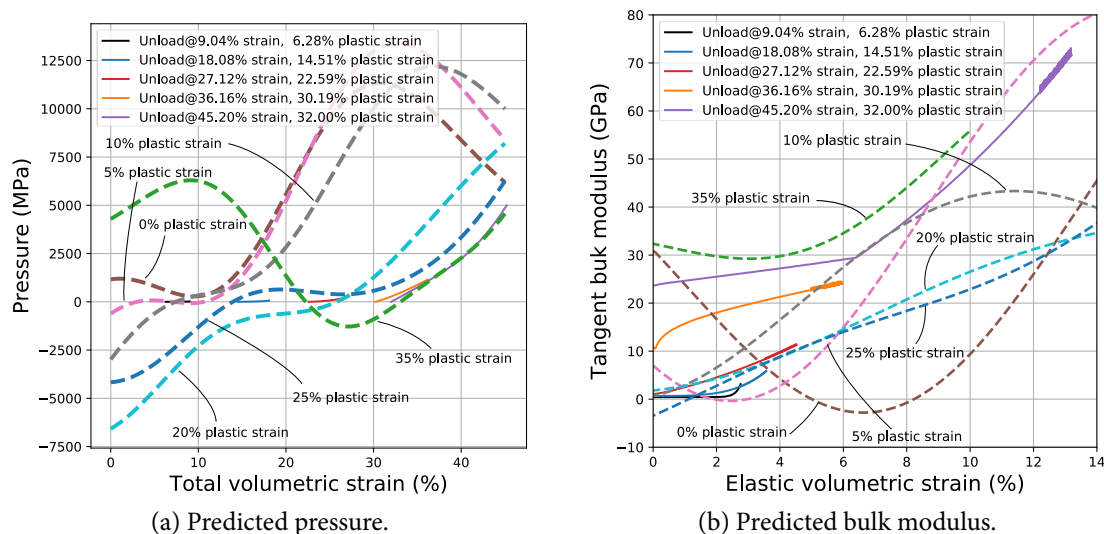


Figure 13 – Support vector regression predictions from fits to total strain data for poorly-graded dry sand.

Predictions from the ν -SVR algorithm (Schölkopf et al., 2000) are shown in Figures 14(a) and (b). Though the response is better behaved, negative initial bulk moduli are predicted for plastic volumetric strains less than approximately 6%. Also, the pressures predicted for the various plastic volumetric strains are nonphysical.

5.3 Extended data, overfitting and cross-validation

Since the strain regime of interest in sand simulations is predominantly compressive, we require the SVR fitted model to be accurate and predictive in that regime. The results from the previous sections show that even though SVR fits to the training data are reasonable to a large extent, predictions outside the training set are inadequate for modeling the material.

We can attempt to obtain better predictions by extending the input data in the tension regime by linearly extrapolating the inputs as depicted by dashed lines in Figure 15. Support vectors are fit to the extended data after shuffling but without bootstrapping.⁵

In the previous sections, the entire input data set was used to fit (train) SVR models. Though the fits to the input data were excellent, the models failed to generalize adequately to plastic strains not included in the training set. This problem is common in machine learning and is typically addressed

⁵If bootstrapping is used on the extended data, the number of samples can become larger than 10,000. Because the LIBSVM implementation is serial, bootstrapped data containing 10,000 or more samples take several hours to fit on a single core of an Intel i7 processor without producing any significant improvements to the fits and predictions.

DRAFT

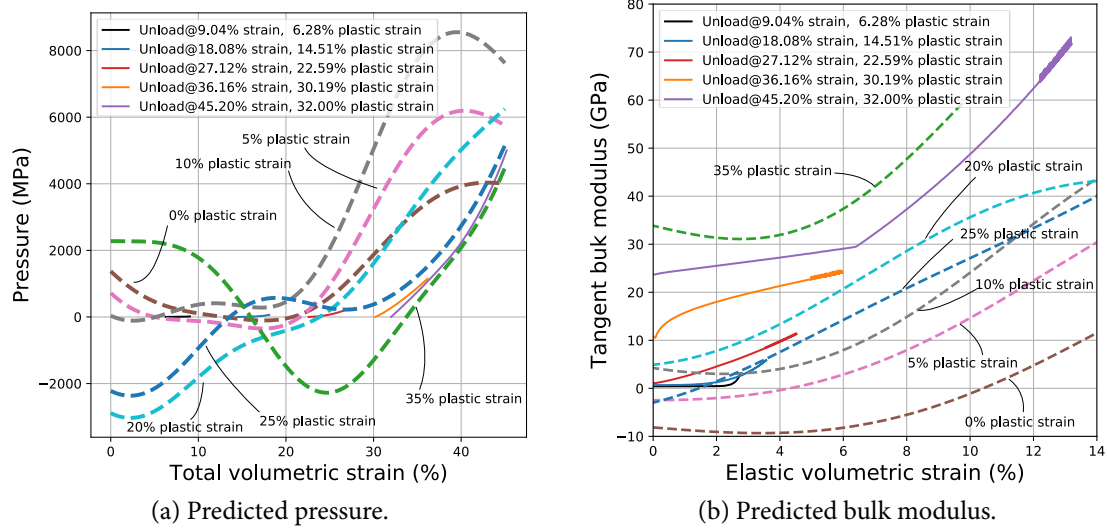


Figure 14 – *v*-SVR (Schölkopf et al., 2000) predictions from fits to total strain data for poorly-graded dry sand.

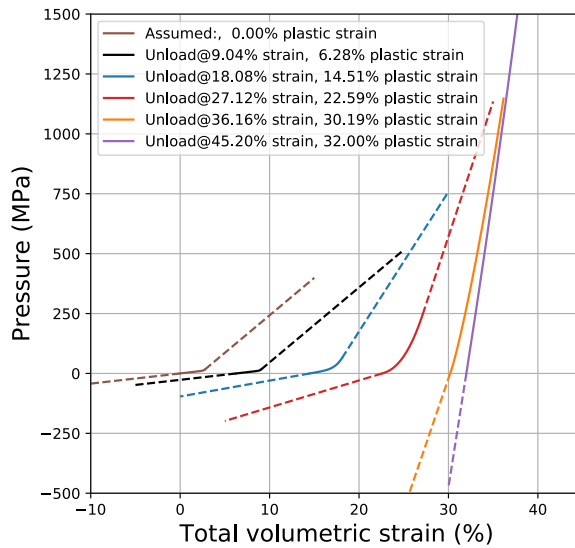


Figure 15 – Extended input data (dashed lines) using linear extrapolation.

by splitting the input data into test and training sets and using cross-validation to test the quality of generalization.

Since the amount of data available is small, we cannot keep aside, during training, all the data for a given plastic strain so that they can later be used to test the quality of the fit. Instead, we keep aside 40% of the combined data set for testing and train the model on the remainder. Shuffled cross-validation on several randomly chosen training data sets (ten sets in our tests) provide an estimate of the quality of generalizability for the plastic strains in the input data. An optimal value of ϵ can be extracted from the process. However, cross-validation is inadequate if we wish to generate a model that is generalizable to plastic strains that are not included in the training set and other metrics are required to select the best model as we shall see later in this section.

Compare the ϵ -SVR fits shown in Figures 16(a) and (b) which have been trained with $C = 1$ and $\epsilon = 0.001$ on 60% of the input data with the plots in Figures 16(c) and (d) that used $C = 10$. The SVR model fits the data better at $C = 10$, though the fits are worse than those seen in the previous section where the full data set was used for training. Notably, no negative bulk moduli are predicted and monotonic increase with strain is observed. The observations suggest that tension-extended data are required for fitting physically reasonable models in the absence of explicit extra constraints in the SVR optimization problem.

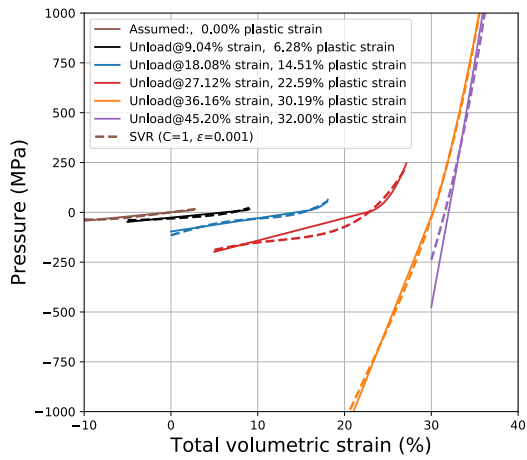
When pressure vs. total volumetric strain curves are computed using these fitted models for plastic strains that are not in the training/test set, we get the predicted curves shown in Figure 17(a) and (b). Visual examination of the predicted curves at 5% and 10% plastic strain suggests that the predicted pressures for $C = 1$ are a better generalization of the input data than those for $C = 10$, even though the fits to the training data are better for $C = 10$. Models that fit the training data better can be found for higher values of C combined with smaller values of ϵ , and a grid search indicated that the least error was for $C = 10,000$ and $\epsilon = 0.0001$. However, for those models the generalization to plastic sets outside the training values was poorer.

Figure 18(a) and (b) show the corresponding bulk modulus predictions for various plastic strains. In this case, the generalization is marginally better for $C = 10$. Given that the fit to the training data is also better at this value of C , we can choose this SVR model as the optimal one given the input data. Though the quality of the fits to the input data from this model are not the best possible, the clear improvement in generalization capability makes it the best choice under the circumstances.

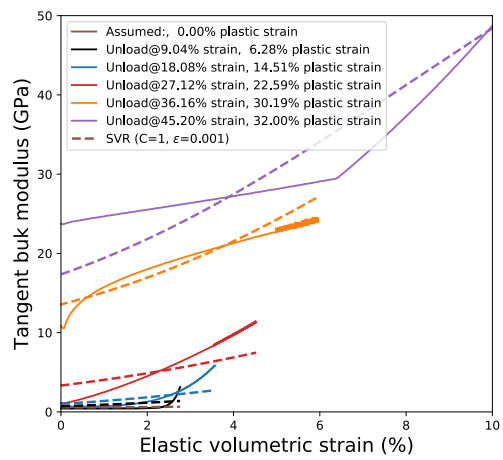
6 Concluding remarks

Differences in material properties that appear small can have a significant effect on the mechanical response of materials, particularly in the non-linear and high strain-rate regimes. We have illustrated the effect with a simple rectangular punch impacting sand. As a consequence, we expect predictive simulations to be sensitive to the accuracy of material models, at least for the case we have discussed in this report. Nonlinear material models are often purely phenomenological, in which case the development of an accurate model involves two steps. In the first step, the model developer explores various functional forms that may possibly model the data and chooses one that appears suitable (and also satisfies constraints imposed by the second law of thermodynamics). Next, model parameters are determined via a convex optimization step. The model function depends strongly on the particular set of data being examined. It is not unusual to see both steps being repeated for what

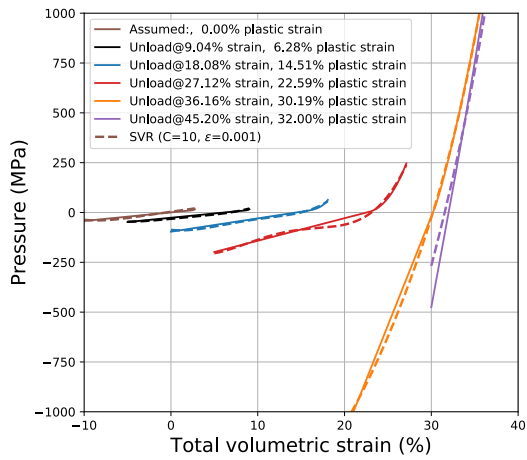
DRAFT



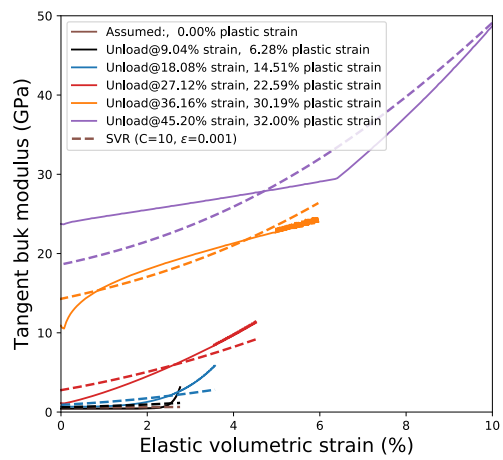
(a) Pressure ($C = 1, \epsilon = 0.001$).



(b) Bulk modulus ($C = 1, \epsilon = 0.001$).

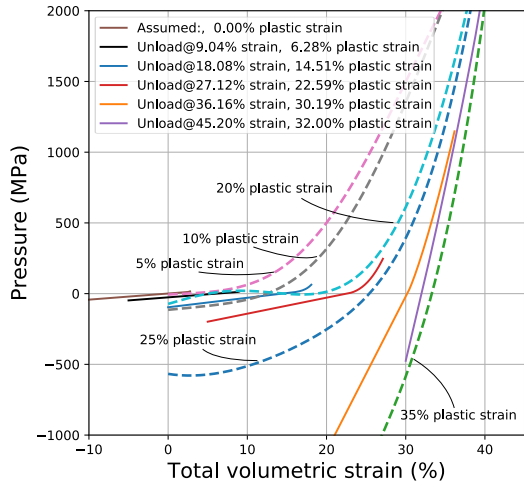


(c) Pressure ($C = 10, \epsilon = 0.001$).

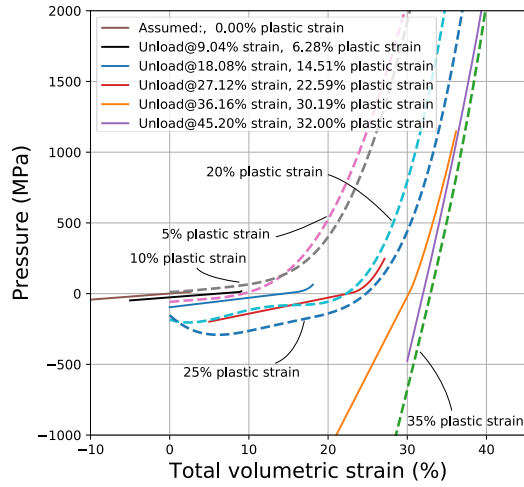


(d) Bulk modulus ($C = 10, \epsilon = 0.001$).

Figure 16 – ϵ -SVR fits to tension extended elastic unloading curves for poorly-graded dry sand.

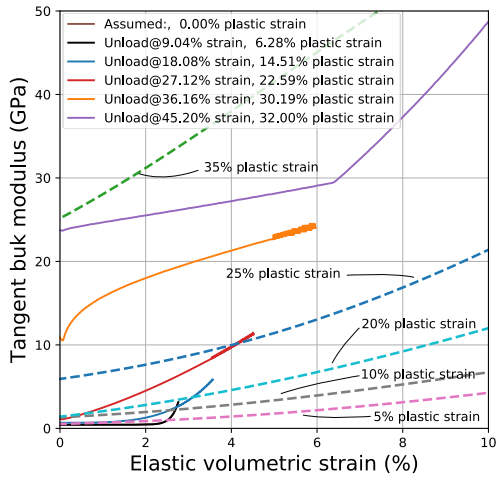


(a) $C = 1, \epsilon = 0.001$.

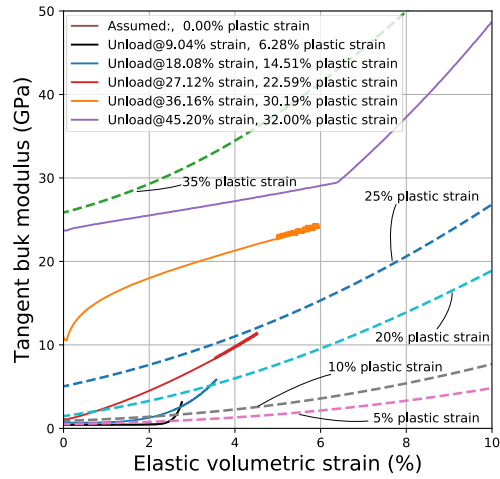


(b) $C = 10, \epsilon = 0.001$.

Figure 17 – ϵ -SVR prediction of pressure vs. elastic strain at various values of plastic strain outside the training set for poorly-graded dry sand.



(a) $C = 1, \epsilon = 0.001$.



(b) $C = 10, \epsilon = 0.001$.

Figure 18 – ϵ -SVR prediction of bulk modulus vs. elastic strain at various values of plastic strain outside the training set for poorly-graded dry sand.

appears nominally to be the same material. This is particularly true of granular materials.

In this report, we suggest that the two-stage model development process may be reduced to one step. In particular, we explore the support vector regression approach. Future work will present similar studies on multi-layer perceptron neural networks. Support-vector regression requires just an optimization step because the functional form is fixed at the outset. However, as we have seen, fitting the model accurately to the input data is not a guarantee that the model will be unusable in a predictive simulation. The issue of overfitting can be addressed by accepting poorer fits to the input data. The support vector model typically generalizes better under those conditions as we demonstrate in this report. The machine learning/statistics literature discusses reasons extensively. However, the selection of a good model still remains somewhat of an art and depends on the scaling of the input data, the choice of support vector kernel, and the judgement of the modeler.

This report has explored a crush-curve model and a nonlinear, plastic strain-dependent bulk modulus model. As we have seen, the bulk modulus model we consider acceptable still fails to fit the input data accurately. Though the model is an improvement over a constant bulk modulus model (or even the Arenisca model), the strong sensitivity of the bulk modulus on sand impact simulations suggests that the support vector model would also be inadequate for predictive simulations. It is unlikely that any purely phenomenological model will be accurate enough and these models should be supplemented by physical models based on micromechanical considerations. Recent progress in discrete element and micromorphic modeling, in conjunction with tabular interpolation, indicates a potential way forward.

Acknowledgements

This research has been partially funded by the US Office of Naval Research PTE Federal award number N00014-17-1-2704.

References

- Adams, Brian M et al. (2009). “DAKOTA, a multilevel parallel object-oriented framework for design optimization, parameter estimation, uncertainty quantification, and sensitivity analysis: version 5.0 user’s manual”. In: *Sandia National Laboratories, Tech. Rep. SAND2010-2183* (cit. on p. 2).
- Banerjee, B. and R. M. Brannon (2017). *Theory, verification, and validation of the ARENA constitutive model for applications to high-rate loading of fully or partially saturated granular media*. Tech. rep. PAR-10021867-1516.v1. Parresia Research Limited and University of Utah. DOI: [10.13140/RG.2.2.10671.53922](https://doi.org/10.13140/RG.2.2.10671.53922) (cit. on pp. 2, 3).
- (2019). “Continuum modeling of partially saturated soils”. In: *Shock Phenomena in Granular and Porous Materials*. Ed. by T. Vogler and A. Fredenburg. Springer (cit. on pp. 2, 3).
- Brannon, R. M. et al. (2015). “KAYENTA: Theory and User’s Guide”. In: *Sandia National Laboratories report SAND2015-0803* (cit. on pp. 1, 8).
- Chang, Chih-Chung and Chih-Jen Lin (2011). “LIBSVM: A library for support vector machines”. In: *ACM transactions on intelligent systems and technology (TIST)* 2.3, p. 27 (cit. on pp. 6, 8).
- Cortes, Corinna and Vladimir Vapnik (1995). “Support-vector networks”. In: *Machine learning* 20.3, pp. 273–297 (cit. on p. 2).

- Fox, D. M. et al. (2014). “The effects of air filled voids and water content on the momentum transferred from a shallow buried explosive to a rigid target”. In: *International Journal of Impact Engineering* 69, pp. 182–193 (cit. on p. 6).
- Hemel, M. A., J. E. Guilkey, and R. M. Brannon (2015). “Continuum effective-stress approach for high-rate plastic deformation of fluid-saturated geomaterials with application to shaped-charge jet penetration”. In: *Acta Mechanica* 227.2, pp. 279–310 (cit. on p. 3).
- (2016). “Continuum effective-stress approach for high-rate plastic deformation of fluid-saturated geomaterials with application to shaped-charge jet penetration”. In: *Acta Mechanica* 227.2, pp. 279–310 (cit. on p. 2).
- Kohestani, VR and M Hassanlourad (2016). “Modeling the mechanical behavior of carbonate sands using artificial neural networks and support vector machines”. In: *International Journal of Geomechanics* 16.1, p. 04015038 (cit. on p. 2).
- McKinney, Wes (2011). “pandas: a foundational Python library for data analysis and statistics”. In: *Python for High Performance and Scientific Computing* 14 (cit. on p. 8).
- Mooney, Christopher Z, Robert D Duval, and Robert Duvall (1993). *Bootstrapping: A nonparametric approach to statistical inference*. 94-95. Sage (cit. on p. 10).
- Pedregosa, F. et al. (2011). “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12, pp. 2825–2830 (cit. on p. 8).
- Schölkopf, Bernhard et al. (2000). “New support vector algorithms”. In: *Neural computation* 12.5, pp. 1207–1245 (cit. on pp. 3, 6, 13, 15, 16).
- Smola, Alex J and Bernhard Schölkopf (2004). “A tutorial on support vector regression”. In: *Statistics and computing* 14.3, pp. 199–222 (cit. on p. 3).
- Vapnik, Vladimir (1998). “The support vector method of function estimation”. In: *Nonlinear Modeling*. Springer, pp. 55–85 (cit. on p. 5).
- (2013). *The nature of statistical learning theory*. Springer science & business media (cit. on p. 2).
- Xue, Long et al. (2016). “In situ identification of shearing parameters for loose lunar soil using least squares support vector machine”. In: *Aerospace Science and Technology* 53, pp. 154–161 (cit. on p. 2).
- Yeo, In-Kwon and Richard A Johnson (2000). “A new family of power transformations to improve normality or symmetry”. In: *Biometrika* 87.4, pp. 954–959 (cit. on p. 8).
- Yuvaraj, P et al. (2013). “Support vector regression based models to predict fracture characteristics of high strength and ultra high strength concrete beams”. In: *Engineering fracture mechanics* 98, pp. 29–43 (cit. on p. 2).
- Zhang, Limao et al. (2017). “Intelligent approach to estimation of tunnel-induced ground settlement using wavelet packet and support vector machines”. In: *Journal of Computing in Civil Engineering* 31.2, p. 04016053 (cit. on p. 2).
- Zhao, Hong-bo and Shunde Yin (2009). “Geomechanical parameters identification by particle swarm optimization and support vector machine”. In: *Applied Mathematical Modelling* 33.10, pp. 3997–4012 (cit. on p. 2).