



A Comparative Study of Feature Selection Methods for Stress Hotspot Classification in Materials

Ankita Mangal¹ · Elizabeth A. Holm¹ 

Received: 19 April 2018 / Accepted: 16 May 2018 / Published online: 15 June 2018
© The Minerals, Metals & Materials Society 2018, corrected publication August/2018

Abstract

The first step in constructing a machine learning model is defining the features of the dataset that can be used for optimal learning. In this work, we discuss feature selection methods, which can be used to build better models, as well as achieve model interpretability. We applied these methods in the context of stress hotspot classification problem, to determine what microstructural characteristics can cause stress to build up in certain grains during uniaxial tensile deformation. The results show how some feature selection techniques are biased and demonstrate a preferred technique to get feature rankings for physical interpretations.

Keywords Stress hotspots · Machine learning · Random forests · Crystal plasticity · Titanium alloys · Feature selection

Introduction

Statistical learning methods are gaining popularity in the materials science field, rapidly becoming known as “Materials Data Science.” With new data infrastructure platforms like Citrination [1] and the Materials data curation system [2], machine learning (ML) methods are entering the mainstream of materials science. Materials data science and informatics is an emergent field aligned with the goals of the Materials Genome Initiative to reduce the cost and time for materials design, development, and deployment. Building and interpreting machine learning models are indispensable parts of the process of curating materials knowledge. ML methods have been used for predicting a target property such as material failure [3, 4], twinning deformation [5], phase diagrams [6], and guiding experiments and calculations in composition space [7, 8]. Machine learning models are built on learning from “features” or variables

that describe the problem. Thus, an important aspect of the machine learning process is to determine which variables most enable data-driven insights about the problem.

Dimensionality reduction techniques (such as principal component analysis (PCA) [9], kernel PCA [10], autoencoders [11], feature compression from information gain theory [12]) have become popular for producing compact feature representations [13]. They are applied to the feature set to get the best feature representation, resulting in a smaller dataset, which speeds up the model construction [14]. Dimensionality reduction has been used by material scientists to establish process-structure-property relationships and for exploratory data analysis to understand trends in a multivariate space [15]. For example, ranking-based feature selection methods such as information gain and Pearson correlation have been used during construction of predictive models for fatigue strength of steel [16]. Kalidindi et al. [17] have used two-point correlations and PCA to describe microstructure-property relationships between local neighborhoods and the localizations in microstructural response. Dey et al. [18] used PCA to analyze the features that cause outliers when predicting bandgaps for new chalcopyrite compounds. Broderick et al. [19] demonstrate how a compact representation (via PCA) makes it easy to visually track the different chemical processing pathways for interpenetrating polymer networks (IPNs) due to changing composition versus changing polymerization. However, dimensionality reduction techniques change the original representation of the features, and hence

The original version of this article was revised: The formatting of Table 2, column 1 in the original version did not match the descriptions in the text and table footer.

✉ Ankita Mangal
mangalanks@gmail.com

Elizabeth A. Holm
eaholm@andrew.cmu.edu

¹ Department of Materials Science and Engineering,
Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh,
PA 15213, USA

offer limited interpretability [13]. An alternate method for better models is feature selection. Feature selection is the process of selecting a subset of the original variables such that a model built on data containing only these features has the best performance. Feature selection avoids overfitting, improves model performance by getting rid of redundant features, and has the added advantage of keeping the original feature representation, thus offering better interpretability [13].

Feature selection methods have been used extensively in the field of bioinformatics [20], psychiatry [21], and cheminformatics [22]. There are multiple feature selection methods, broadly categorized into filter, wrapper, and embedded methods based on their interaction with the predictor during the selection process. The filter methods rank the variables as a preprocessing step, and feature selection is done before choosing the model. In the wrapper approach, nested subsets of variables are tested to select the optimal subset that work best for the model during the learning process. Embedded methods are those which incorporate variable selection in the training algorithm.

We have used random forest models to study stress hotspot classification in FCC [3] and HCP [4] materials. In this paper, we review some feature selection techniques applied to the stress hotspot prediction problem in hexagonal close-packed materials, and compare them with respect to future data prediction. We focus on two commonly used techniques from each method: (1) filter methods such as correlation-based feature selection (CFS) [23] and Pearson correlation [24]; (2) wrapper methods such as FeaLect [25] and recursive feature elimination (RFE) [13]; and (3) embedded methods such as random forest permutation accuracy importance (RF-PAI) [26] and least absolute shrinkage and selection operator (LASSO) [27]. The main contribution of this article is to raise awareness in the materials data science community about how different feature selection techniques can lead to misguided model interpretations and how to avoid them. We point out some of the inadequacies of popular feature selection methods, and finally, we extract data-driven insights with better understanding of the methods used.

Methods

An applied stress is distributed heterogeneously within the grains in a microstructure [28]. Under an applied deformation, some grains are prone to accumulating stress due to their orientation, geometry, and placement with respect to the neighboring grains. These regions of high stress, so-called stress hotspots, are related to void nucleation under ductile fracture [29]. Stress hotspot formation has been studied in face-centered cubic (FCC)

[3] and hexagonal close-packed (HCP) [4] materials using a machine learning approach. A set of microstructural descriptors was designed to be used as features in a random forest model for predicting stress hotspots. To achieve data-driven insights into the problem, it is essential to rank the microstructural descriptors (features). In this paper, we review different feature selection techniques applied to the stress hotspot classification problem in HCP materials, which have a complex plasticity landscape due to anisotropic slip system activity.

Let (x_i, y_i) , for $i = 1, \dots, N$ be N independent identically distributed (i.i.d.) observations of a p -dimensional vector of grain features $x_i \in R^p$, and the response variable $y_i \in \{0, 1\}$ denotes the truth value of a grain being a stress hotspot. The input matrix is denoted by $X = (x_1, \dots, x_N) \in R^{N \times p}$, and $y \in \{0, 1\}^N$ is the binary outcome. We will use small letters to refer to the samples x_1, \dots, x_N and capital letters to refer to the features X_1, \dots, X_p of the input matrix X . Feature importance refers to metrics used by various feature selection methods to rank, such as feature weights in linear models or variable importance in random forest models.

Dataset Studied

The machine learning input dataset of synthetic 3D equiaxed microstructures with different HCP textures was generated using Dream.3D in [4]. Uniaxial tensile deformation was simulated in these microstructures using EVPFFT [30] with constitutive parameters representing a titanium-like HCP material with an anisotropic critically resolved shear stress ratio [4]. The EVPFFT simulation was carried out in 200 strain steps of 0.01% along sample Z direction, up to a total strain of 2%. The crystal plasticity simulations result in spatially resolved micromechanical stress and strain fields. This data was averaged to attain a dataset containing grain-wise values for equivalent Von Mises stress, and the corresponding Euler angles and grain connectivity parameters. Steps to reproduce this dataset are discussed in detail in [31].

The grains having stress greater than the 90th percentile of the stress distribution within each microstructure are designated as stress hotspots, a binary target. Thirty-four variables to be used as features in machine learning were developed. These features (X) describe the grain texture and geometry and have been summarized in Table 1. We note that these features are not a complete set, and there are long-range effects causing stress hotspots. We have taken the first-order microstructural descriptors to build stress hotspot prediction models and understand that these models can be improved upon by adding the missing features.

The microstructures contained in this dataset represent eight different kinds of textures, and we validate the machine learning models by leaving one texture out

Table 1 Feature name descriptions

Feature name abbreviation	Description	Feature name abbreviation	Description
Schmid_1	Basal $\langle a \rangle$ Schmid factor	100_IPF_x	Distance of tensile axis from the corners of the 100 Inverse pole figure
Schmid_2	Prismatic $\langle a \rangle$ Schmid factor	001_IPF_x	Distance of tensile axis from the corners of the 001 Inverse pole figure
Schmid_3	Pyramidal $\langle a \rangle$ Schmid factor	010_IPF_x	Distance of tensile axis from the corners of the 010 Inverse pole figure
Schmid_4	Pyramidal $\langle c + a \rangle$ Schmid factor	AvgC_Axes_x	Unit vector components describing the c axis orientation for hcp
Max_mis	Maximum misorientation between a grain and its nearest neighbor	Min_mis	Minimum misorientation between a grain and its nearest neighbor
AvgMisorientations	Average misorientation between a grain and its nearest neighbor	KernelAvg	Average misorientation within a grain
θ	Polar angle of hcp c axis w.r.t sample frame	ϕ	Azimuthal Angle of hcp c axis w.r.t. sample frame
GBEuc	Average distance of grain voxels to the nearest grain boundary	TJEuc	Average distance of grain voxels to the nearest triple line
QPEuc	Average distance of grain voxels to the nearest quadruple point		
FeatureVolumes	Volume of grain	EquivalentDiameters	Equivalent spherical diameter of a grain
NumNeighbors	Number of nearest neighbors of a grain	Neighborhoods	Number of grains having their centroid within the 1 multiple of equivalent sphere diameters from each grain

validation. This divides the dataset into $\sim 85\%$ training and $\sim 15\%$ validation. Note that since only 10% of the grains are stress hotspots, this is an imbalanced classification problem. Hence, the model performance is measured by the AUC (area under curve), a metric for binary classification which is insensitive to imbalance in the classes. An AUC of 100% denotes perfect classification and 50% denotes no better than random guessing [32].

We first build a decision tree-based random forest model [26] for stress hotspot classification using all the thirty-four variables. We then rank and select the variables using different feature selection techniques. The selected variables are then used to build random forest models and we observe the improvement in model performance. The feature rankings are then used to gain insights about the physics behind stress hotspot formation.

Feature Selection Methods

Filter Methods

Filter methods are based on preprocessing the dataset to extract the features X_1, \dots, X_p that most impact the target Y . Some of these methods are as follows:

Pearson Correlation [24] This method provides a straightforward way for filtering features according to their correlation coefficient. The Pearson correlation coefficient between a feature X_i and the target Y is:

$$\rho_i = \frac{\text{cov}(X_i, Y)}{\sigma(X_i)\sigma_Y}$$

where $\text{cov}(X_i, Y)$ is the covariance, σ is the standard deviation [24]. It ranges between $(-1, 1)$ from negative to positive correlation, and can be used for binary classification and regression problems. It is a quick metric using which the features are ranked in order of the absolute correlation coefficient to the target.

Correlation-based feature selection (CFS) [23] CFS was developed to select a subset of features with high correlation to the target and low intercorrelation among themselves, thus reducing redundancy and selecting a diverse feature set. CFS gives a heuristic merit over a feature subset instead of individual features. It uses symmetrical uncertainty correlation coefficient given by:

$$r(X, Y) = 2.0 \times \frac{IG(X|Y)}{H(X) + H(Y)}$$

where $IG(X|Y)$ is the information gain of feature X for the class attribute Y . $H(X)$ is the entropy of variable X . The following merit metric was used to rank each subset S containing k features:

$$Merits_S = \frac{k\bar{r}_{cf}}{\sqrt{k + k(k-1)\bar{r}_{ff}}}$$

where \bar{r}_{cf} is the mean symmetrical uncertainty correlation between the feature ($f \in S$) and the target, and \bar{r}_{ff} is the average feature-feature intercorrelation. To account for the high computational complexity of evaluating all possible feature subsets, CFS is often combined with search strategies such as forward selection, backward elimination, and bidirectional search. In this work, we have used the scikit-learn implementation of CFS [33] which uses symmetrical uncertainty [23] as the correlation metric and explores the subset space using best first search [34], stopping when it encounters five consecutive fully expanded non-improving subsets.

Embedded Methods

These methods are popular because they perform feature selection while constructing the classifier, removing the preprocessing feature selection step. Some popular algorithms are support vector machines (SVM) using RFE [35], RF [26], and LASSO [27]. We compare LASSO and RF methods for feature selection on the stress hotspot dataset.

Least Absolute Shrinkage and Selection Operator (LASSO)

[27] LASSO is linear regression with L_1 regularization [27]. A linear model \mathcal{L} is constructed

$$\mathcal{L} : \min_{w \in \mathbb{R}^p} \sum_{i=1}^N \frac{1}{2N} \|y_i - w^T \cdot x_i\|_2^2 + \lambda \|w\|_1$$

on the training data (x_i, y_i) , $i = 1, \dots, N$, where w is a p dimensional vector of weights corresponding to each feature dimension p . The L_1 regularization term ($\lambda \|w\|_1$) helps in feature selection by pushing the weights of correlated features to zero, thus preventing overfitting and improving model performance. Model interpretation is possible by ranking the features according to the LASSO feature weights. However, it has been shown that for a given regularization strength λ , if the features have redundancy, inconsistent subsets can be selected [36]. Nonetheless, LASSO has been shown to provide good prediction accuracy by reducing model variance without substantially increasing the bias while providing better model interpretability. We used the scikit-learn implementation to compute our results [37].

Random Forest Permutation Accuracy Importance (RF-PAI)

[26] The random forest is a nonlinear multivariate model

built on an ensemble of decision trees. It can be used to determine feature importance using the inbuilt feature importance measure [26]. For each of the trees in the model, a feature node is randomly replaced with another feature node while keeping all other nodes unchanged. The resulting model will have a lower performance if the feature is important. When the permuted variable X_j , together with the remaining unchanged variables, is used to predict the response, the number of observations classified correctly decreases substantially, if the original variable X_j was associated with the response. Thus, a reasonable measure for feature importance is the difference in prediction accuracy before and after permuting X_j . The feature importance calculated this way is known as PAI and was computed using the scikit-learn package in Python [37].

Wrapper Methods

Wrapper methods test feature subsets using a model hypothesis. Wrapper methods can detect feature dependencies, i.e., features that become important in the presence of each other. They are computationally expensive, hence often used in greedy search strategies (forward selection and backward elimination [38]) which are fast and avoid overfitting to get the best nested subset of features.

FeaLect Algorithm [25] The number of features selected by LASSO depends on the regularization parameter λ , and in the presence of highly correlated features, LASSO arbitrarily selects one feature from a group of correlated features [39]. The set of possible solutions for all LASSO regularization strengths is given by the regularization path, which can be recovered computationally efficiently using the least angles regression (LARS) algorithm [40]. It was shown that LASSO selects the relevant variables with a probability 1 and all other with a positive probability [36]. An improvement in LASSO, the Bolasso feature selection algorithm was developed based on this property [36] in 2008. In this method, the dataset is bootstrapped, and a LASSO model with a fixed regularization strength λ is fit to each subset. Finally, the intersection of the LASSO-selected features in each subset is chosen to get a consistent feature subset.

In 2013, the FeaLect algorithm, an improvement over the Bolasso algorithm, was developed based on the combinatorial analysis of regression coefficients estimated using LARS [25]. FeaLect considers the full regularization path and computes the feature importance using a combinatorial scoring method, as opposed to simply taking the intersection with Bolasso. The FeaLect scoring scheme measures the quality of each feature in each bootstrapped sample and averages them to select the most relevant features,

providing a robust feature selection method. We used the R implementation of FeaLect to compute our results [41].

Recursive Feature Elimination (RFE) [35] A number of common ML techniques (such as linear regression, SVM, decision trees, Naive Bayes, perceptron, etc.) provide feature weights that consider multivariate interacting effects between features [13]. To interpret the relative importance of the variables from these model feature weights, RFE was introduced in the context of SVM [35] for getting compact gene subsets from DNA-microarray data.

To find the best feature subset, instead of doing an exhaustive search over all feature combinations, RFE uses a greedy approach, which has been shown to reduce the effect of correlation bias in variable importance measures [42]. RFE uses backward elimination by taking the given model (SVM, random forests, linear regression, etc.) and discarding the worst feature (by absolute classifier weight or feature ranking), and repeating the process over increasingly smaller feature subsets until the best model hypothesis is achieved. The weights of this optimal model are used to rank features. Although this feature ranking might not be the optimal ranking for individual features, it is often used as a variable importance measure [42]. We used the scikit-learn implementation of RFE with random forest classifier to come up with a feature ranking for our dataset.

Results and Discussion

Table 2 shows the feature importances calculated using filter-based methods Pearson correlation and CFS; embedded methods RF, linear regression, ridge regression (L_2 regularization), and LASSO regression; and finally wrapper methods RFE and FeaLect. The values in bold font denote the features that were finally selected to build RF models and their corresponding performances are noted. The input data was scaled by minimum and maximum values to [0,1]. Figure 1 shows the correlation matrix for the features and the target.

Pearson correlation can be used for feature selection, resulting in a good model. However, this measure has implicit orthogonality assumptions between variables, and the coefficient does not take mutual information between features into account. Additionally, this method only looks for linear correlations which might not capture many physical phenomenon.

The feature subset selected by CFS contains features with higher class correlation and lower redundancy, which translate to a good predictive model. Although we know grain geometry and neighborhood are important to hotspot formation, CFS does not select any geometry-based features and fails to provide an individual feature ranking.

Linear regression, ridge regression, and LASSO are highly correlated linear models. A simple linear model results in huge weights for some features (NumCells, FeatureVolumes), likely due to overfitting, and hence is unsuitable for deducing variable importance. Ridge regression compensates for this problem by using L_1 regularization, but the weights are distributed among the redundant features, which might lead to incorrect conclusions. LASSO regression overcomes this problem by pushing the weights of correlated features to 0, resulting in a good feature subset. The top five ranked features by LASSO with regularization strength of $\lambda = 0.3$ are : $\sin\theta$, AvgMisorientations, $\cos\phi$, $\sin\phi$ and Schmid_1. The first geometry-based feature ranks 10th on the list, which seems to underestimate the physical importance of such features. A drawback of deriving insights from LASSO-selected features is that it arbitrarily selects a few representatives from the correlated features, and the number of features selected depends heavily on the regularization strength. Thus, the models become unstable, because changes in training subset can result in different selected features. Hence, these methods are not ideal for deriving physical insights from the model.

Random forest models also provide an embedded feature ranking module. The RF-PAI importance seems to focus only on the hcp “c” axis orientation-derived features ($\cos\phi$, $\sin\theta$), average misorientation, and the Prismatic $< a >$ Schmid factor, while discounting most of the geometry-derived features. RF-PAI suffers from correlation bias due to preferential selection of correlated features during tree building process [43]. As the number of correlated variables increases, the feature importance score for each variable decreases. Oftentimes, the less relevant variables replace the predictive ones (due to correlation) and thus receive undeserved, boosted importance [44]. Random forest variable importance can also be biased in situations where the features vary in their scale of measurement or number of categories, because the underlying Gini gain splitting criterion is a biased estimator and can be affected by multiple testing effects [45]. From Fig. 1, we found that all the geometry-based features are highly correlated to each other, therefore deducing physical insights from this ranking is unsuitable.

Hence, we move to wrapper-based methods for feature importance. RFE has been shown to reduce the effect of the correlation on the importance measure [42]. RFE with underlying random forest model selects a feature subset consisting of two geometry-based features (GBEuc and EquivalentDiameter); however, it fails to give an individual ranking among the features.

FeaLect provides a robust feature selection method by compensating for the uncertainty in LASSO due to arbitrary selection among correlated variables, and the number of

Table 2 Variable Importance Measures using different methods for HCP materials with Unequal CRSS

Features	Pearson correlation	CFS	RFE	RF ($\times 1e - 2$)	Linear regression	Ridge regression	LASSO regression	Fealect ($\times 1e - 2$)
cos ϕ	-0.29	1	1	53.43	27.37	27.36	26.01	245.0
Schmid_1	-0.39	0	1	0.15	22.72	22.69	14.78	145.00
EquivalentDiameters*	-0.01	0	1	0.05	0.15	0.15	0.08	79.47
GBEuc*	-0.01	0	1	0.12	0.22	0.22	0.12	71.47
Schmid_4	-0.18	0	1	0.31	7.29	7.31	10.35	41.27
Neighborhoods*	-0.01	0	22	0.01	0.10	0.10	0.00	5.53
sin θ	0.48	1	1	8.74	74.78	74.61	52.99	5.00
TJEuc*	-0.01	0	2	0.07	0.97	0.97	0.44	4.93
sin ϕ	0.14	1	16	0.03	80.46	79.96	19.17	1.0
AvgMisorientations	0.31	0	1	8.95	32.08	32.09	32.05	0.83
NumNeighbors*	-0.01	0	23	0.01	0.18	0.17	0.03	0.50
Schmid_3	0.12	0	9	0.03	4.05	4.04	0.00	0.0
Min_mis	0.09	0	1	0.72	3.46	3.46	2.19	0.0
AvgC_Axes_1	0.00	0	1	0.22	0.09	0.09	0.00	0.0
Max_mis	0.17	0	4	0.02	0.86	0.86	0.03	0.0
NumCells*	-0.01	0	18	0.04	1.3e6	0.11	0.21	0.0
Schmid_2	0.49	0	1	26.80	38.03	37.83	8.37	0.0
KernelAvg	-0.01	0	25	0.0	0.22	0.22	0.00	0.0
010_IPF_1	-0.07	0	5	0.01	0.49	0.49	0.00	0.0
ϕ	0.13	1	3	3.4	66.42	65.94	7.68	0.0
001_IPF_0	0.00	0	11	0.03	0.58	0.57	0.00	0.0
001_IPF_2	0.09	0	21	0.01	0.21	0.24	0.19	0.0
010_IPF_0	0.00	0	12	0.01	0.76	0.76	0.23	0.0
100_IPF_0	0.00	0	10	0.01	0.13	0.13	0.00	0.0
001_IPF_1	0.16	0	15	0.01	0.17	0.14	0	0.0
100_IPF_1	0.07	0	14	0.02	1.10	1.10	0.00	0.0
QPEuc*	-0.01	0	6	0.02	0.57	0.57	0.00	0.0
AvgC_Axes_0	0.00	0	7	0.03	0.34	0.34	0.05	0.0
θ	0.00	1	24	0.02	0.04	0.04	0.00	0.0
FeatureVolumes*	-0.01	0	13	0.04	1.3e6	0.11	0.00	0.0
010_IPF_2	-0.04	0	17	0.01	0.79	0.79	0.00	0.0
AvgC_Axes_2	0.00	0	8	0.01	2.9e4	0.07	0.00	0.0
100_IPF_2	0.04	0	19	0.01	1.21	1.20	0.00	0.0
cos θ	0.00	1	20	0.01	2.9e4	0.07	0.00	0.0

Random Forest model validation AUC without feature selection: 71.94%

Random Forest model AUC with selected features (%)

Training	84.02	82.51	84.24	83.82	84.20	84.19	84.31	84.28
Validation	80.46	80.45	80.73	80.19	80.72	80.61	80.83	80.75

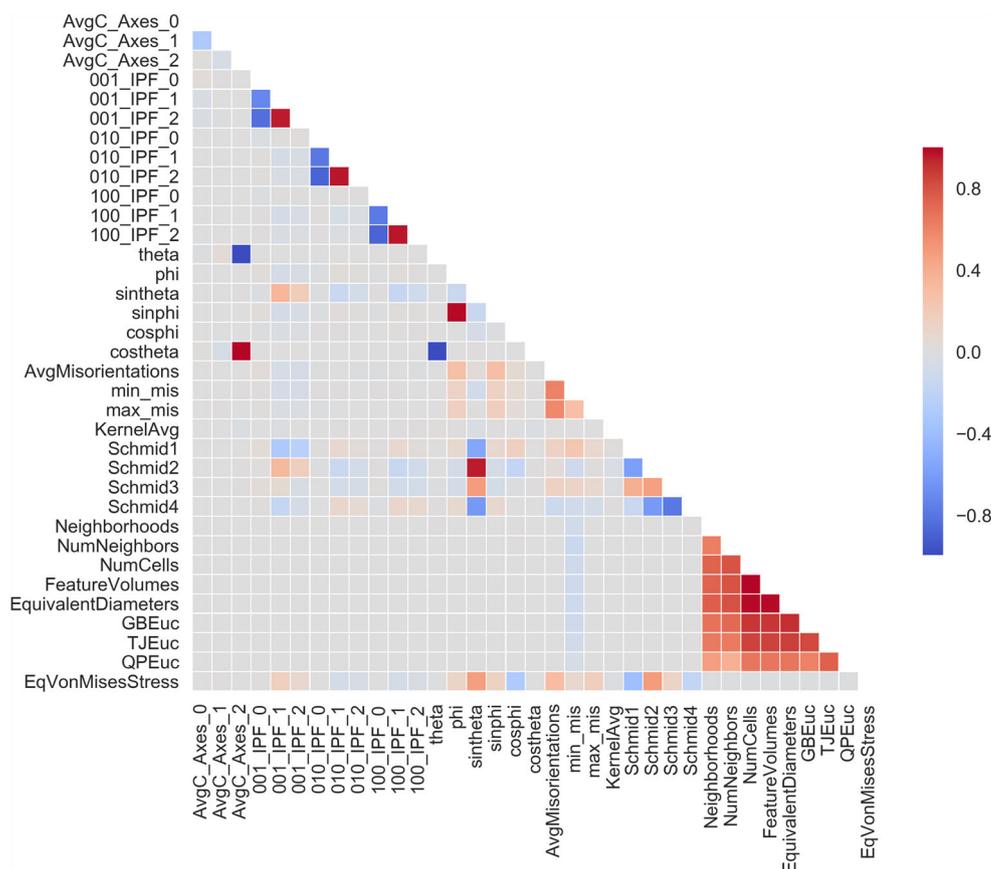
*Feature that describes grain geometry

Values in bold face type indicate features selected to build the Random Forest model by the given technique

selected variables due to change in regularization strength. Table 2 lists the FeaLect-selected variables in decreasing order. We find that the top two important features are derived from the grain crystallography, and geometry-derived features come next. This suggests that both texture- and geometry-based features are important. Using linear

regression-based methods such as these tells us which features are important by themselves, as opposed to RF-PAI which indicates the features that become important due to interactions between them (via RF models) [13]. The FeaLect method provides the best estimate of the feature importance ranking which can then be used to extract

Fig. 1 Pearson correlation matrix between the target (EqVonMisesStress) and all the features



physical insights. This method also divides the features into three classes: informative, irrelevant features that cause model overfitting, and redundant features [25]. The most informative features are the following: $\cos\phi$, $Schmid_1$, EquivalentDiameter, GBEuc, $Schmid_4$, Neighborhoods, $\sin\theta$, and TJEuc. The irrelevant features are $\sin\phi$ and AvgMisorientations (which cause model overfitting). The remaining features are redundant.

A number of selected features directly or indirectly represent the HCP c-axis orientation, such as $\cos\phi$, $\sin\theta$ and basal Schmid factor ($Schmid_1$), which is proportional to $\cos\theta$. It is interesting that pyramidal $\langle c + a \rangle$ Schmid factor ($Schmid_4$) is chosen as important. From Fig. 1, we can see that hot grains form where θ, ϕ maximize $\sin\theta$ and $\sin\phi$, i.e., $\theta \sim 90, \phi \sim 90$. This means that the HCP c-axis orientation of hot grains aligns with the sample Y-axis, which means these grains have a low elastic modulus. Since the c-axis is perpendicular to the tensile axis (sample Z), the deformation along the tensile direction can be accommodated by prismatic slip in these grains, and if pyramidal slip is occurring, it means they have a very high stress [4]. This explains the high importance of the pyramidal $\langle c + a \rangle$ Schmid factor. From the Pearson correlation coefficients in Fig. 1, we can observe that the stress hotspots form in grains with low basal and pyramidal

$\langle c + a \rangle$ Schmid factor, high prismatic $\langle a \rangle$ Schmid factor, and higher values of $\sin\theta$ and $\sin\phi$.

From Fig. 1, we can see that all the grain geometry descriptors do not have a direct correlation with stress, but are still selected by FeaLect. This points to the fact that these variables become important in association with others. We analyzed these features in detail in [4] and found that the hotspots lie closer to grain boundaries (GBEuc), triple junctions (TJEuc), and quadruple points (QPEuc), and prefer to form in smaller grains.

There is a subtle distinction between the physical impact of a variable on the target versus the variables that work best for a given model. From Table 2, we can see that a random forest model built on the entire feature set without feature selection has an AUC of 71.94%. All the feature selection techniques result in an improvement in the performance of the random forest model to a validation AUC of about 81%. However, to draw physical interpretations, it is important to use a feature selection technique which (1) keeps the original representation of the features, (2) is not biased by correlations/redundancies among features, (3) is insensitive to the scale of variable values, (4) is stable to the changes in the training dataset, (5) takes multivariate dependencies between the features into account, and (6) provides an individual feature ranking measure.

Conclusions

In this work, we have surveyed different feature selection techniques by applying them to the stress hotspot classification problem. These techniques can be divided into three categories: filter, embedded, and wrapper. We have explored the most commonly used techniques under each category. It was found that all the techniques lead to an improvement in the model performance and are suitable for feature selection to build a better model. However, when the aim is to interpret the model and understand which features might be more causal than others, it is essential to note the limitations of different techniques. We found that in the presence of correlated features, the FeaLect method helped us determine the underlying importance of the features. We find that:

- All feature selection techniques result in $\sim 9\%$ improvement in the AUC metric for stress hotspot classification.
- Correlation-based feature selection and recursive feature elimination are computationally expensive to run, and give only a feature subset ranking.
- Random forest embedded feature ranking is biased against correlated features and hence should not be used to derive physical insights.
- Linear regression-based feature selection techniques can objectively denote the most important features, however have their flaws. These methods can be affected by the scale of features, correlation between them, and the dataset itself.
- The FeaLect algorithm can compensate for the variability in LASSO regression, providing a robust feature ranking that can be used to derive insights.
- Stress hotspot formation under uniaxial tensile deformation is determined by a combination of crystallographic and geometric microstructural descriptors.
- It is essential to choose a feature selection method that can find this dependence even when features are redundant or correlated.

Acknowledgements This work was performed at Carnegie Mellon University. The authors are grateful to the authors of `skfeature` and `sklearn` python libraries who made their source code available through the Internet. We would also like to thank the reviewers for their thorough work. Ricardo Lebensohn of the Los Alamos National Laboratory is acknowledged for the use of the MASSIF code.

Funding Information This work has been supported by the United States National Science Foundation award number DMR-1307138 and DMR-1507830.

References

1. O'Mara J, Meredig B, Michel K (2016) Materials data infrastructure: A case study of the citrination platform to examine data import, storage, and access. *JOM* 68(8):2031. <https://doi.org/10.1007/s11837-016-1984-0>
2. Dima A, Bhaskarla S, Becker C, Brady M, Campbell C, Dessauw P, Hanisch R, Kattner U, Kroenlein K, Newrock M, Peskin A, Plante R, Li SY, Rigodiat PF, Amaral GS, Trautt Z, Schmitt X, Warren J, Youssef S (2016) Informatics infrastructure for the Materials Genome Initiative. *JOM* 68(8):2053. <https://doi.org/10.1007/s11837-016-2000-4>
3. Mangal A, Holm EA (2018) Applied machine learning to predict stress hotspots I: Face centered cubic materials. arXiv:1711.00118v3
4. Mangal A, Holm EA (2018) Applied machine learning to predict stress hotspots II: Hexagonal close packed materials. arXiv:1804.05924
5. Orme AD, Chelladurai I, Rampton TM, Fullwood DT, Khosravani A, Miles MP, Mishra RK (2016) Insights into twinning in Mg AZ31: A combined EBSD and machine learning study. *Comput Mater Sci* 124:353
6. Ch'Ng K, Carrasquilla J, Melko RG, Khatami E (2017) Machine learning phases of strongly correlated fermions. *Phys Rev X* 7(3):1. <https://doi.org/10.1103/PhysRevX.7.031038>
7. Ling J, Hutchinson M, Antono E, Paradiso S, Meredig B (2017) High-dimensional materials and process optimization using datadriven experimental design with well-calibrated uncertainty estimates. *Integr Mater Manuf Innov* 6(3):207. <https://doi.org/10.1007/s40192-017-0098-z>
8. Oliynyk AO, Antono E, Sparks TD, Ghadbeigi L, Gaultois MW, Meredig B, Mar A (2016) High-throughput machine-learning-driven synthesis of full-Heusler compounds. *Chem Mater* 28(20):7324. <https://doi.org/10.1021/acs.chemmater.6b02724>
9. Wall ME, Rechtsteiner A, Rocha LM (2003) In: A practical approach to microarray data analysis. Springer, Berlin, pp 91–109
10. Mika S, Scholkopf B, Smola A, Muller KR, Scholz M, Riitsch G (1999) In: *Adv. Neural Inf. Process. Syst.*, pp 536–542. <http://papers.nips.cc/paper/1491-kernel-pca-and-de-noising-in-feature-spaces.pdf>
11. Hinton GE, Salakhutdinov RR (2006) Reducing the dimensionality of data with neural networks. *Science* 313(80-):504. <https://doi.org/10.1126/science.1127647>
12. Yu L, Liu H (2003) In: *Proceedings of the 20th International Conference in Machine Learning*, pp 856–863. <https://doi.org/citeulike-article-id:3398512>. <http://www.aai.org/Papers/ICML/2003/ICML03-111.pdf>
13. Guyon I, Elisseeff A (2003) An introduction to variable and feature selection. *J Mach Learn Res* 3(3):1157. <https://doi.org/10.1016/j.aca.2011.07.027>
14. Van Der Maaten L, Postma E, Van Den Herik J (2009) Dimensionality reduction : A comparative review. *J Mach Learn Res* 10(2009):66. <https://doi.org/10.1080/1350628044000102>. <http://www.uvt.nl/ticc>
15. Rajan K, Suh C, Mendez PF (2009) Principal component analysis and dimensional analysis as materials informatics tools to reduce dimensionality in materials science and engineering. *Stat Anal Data Min ASA Data Sci J* 1(6):361. <https://doi.org/10.1002/sam>
16. Agrawal A, Deshpande PD, Cecen A, Basavarsu GP, Choudhary AN, Kalidindi SR (2014) Exploration of data science techniques to predict fatigue strength of steel from composition and processing parameters. *Integr Mater Manuf Innov* 3(8):1. <https://doi.org/10.1186/2193-9772-3-8>
17. Kalidindi SR, Niezgodá SR, Salem AA (2011) Microstructure informatics using higher-order statistics and efficient data-mining protocols. *JOM* 63(4):34–41
18. Dey P, Bible J, Datta S, Broderick S, Jasinski J, Sunkara M, Rajan K (2014) Informatics-aided bandgap engineering for solar materials. *Comput Mater Sci* 83:185–195

19. Broderick SR, Nowers JR, Narasimhan B, Rajan K (2009) Tracking chemical processing pathways in combinatorial polymer libraries via data mining. *J Comb Chem* 12(2):270. <https://doi.org/10.1021/cc900145d>
20. Saeyns Y, Inza I, Larranaga P (2007) Gene expression A review of feature selection techniques in bioinformatics. *Bioinformatics* 23(19):2507. <https://doi.org/10.1093/bioinformatics/btm344>
21. Lu F, Petkova E (2014) A comparative study of variable selection methods in the context of developing psychiatric screening instruments. *Stat Med* 33(3):401. <https://doi.org/10.1002/sim.5937>
22. Wegner JK, Fröhlich H, Zell A (2004) Feature selection for descriptor based classification models. 1. Theory and GA-SEC algorithm. *J Chem Inf Comput Sci* 44(3):921. <https://doi.org/10.1021/ci0342324>
23. Hall MA, Smith LA (1999) Feature selection for machine learning: comparing a correlation-based filter approach to the wrapper. In: FLAIRS conference, vol 1999, pp 235–239. <https://pdfs.semanticscholar.org/31ff/33fadae7b0b3a5608a85a35f84ed74659569.pdf>
24. Cohen I, Huang Y, Chen J, Benesty J (2009) In: Noise reduction in speech processing. Springer, pp 1–4. <https://doi.org/10.1007/978-3-642-00296-0>
25. Zare H, Haffari G, Gupta A, Brinkman RR (2013) Scoring relevancy of features based on combinatorial analysis of Lasso with application to lymphoma diagnosis. *BMC Genom* 14(Suppl 1):S14. <https://doi.org/10.1186/1471-2164-14-S1-S14>. <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3549810&tool=pmcentrez&rendertype=abstract>
26. Breiman L (1996) Out-of-bag-estimation. <https://doi.org/10.1007/s13398-014-0173-7.2>
27. Tibshirani R (1996) Regression selection and shrinkage via the lasso. <https://doi.org/10.2307/2346178>. <http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.35.7574>
28. Qidwai MAS, Lewis AC, Geltmacher AB (2009) Using image-based computational modeling to study microstructure – yield correlations in metals. *Acta Mater* 57(14):4233. <https://doi.org/10.1016/j.actamat.2009.05.021>
29. Hull D, Rimmer DE (1959) The growth of grain-boundary voids under stress. *Philos Mag* 4(42):673. <https://doi.org/10.1080/14786435908243264>
30. Lebensohn RA, Kanjarla AK, Eisenlohr P (2012) An elastoviscoplastic formulation based on fast Fourier transforms for the prediction of micromechanical fields in polycrystalline materials. *Int J Plast* 59:32–33. <https://doi.org/10.1016/j.ijplas.2011.12.005>
31. Mangal A, Holm EA (2018) A dataset of synthetic hexagonal close packed 3D polycrystalline microstructures, grain-wise microstructural descriptors and grain averaged stress fields under uniaxial tensile deformation for two sets of constitutive parameters. (in preparation for Data in Brief)
32. Hanley JA, McNeil BJ (1982) The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143(1):29. <https://doi.org/10.1148/radiology.143.1.7063747>
33. Zhao Z, Morstatter F, Sharma S, Alelyani S, Anand A, Liu H (2010) Advancing Feature Selection Research, ASU Featur. Sel. Repos. Arizona State University, pp 1 – 28. http://featureselection.asu.edu/featureselection_techreport.pdf
34. Pearl J (1984) Heuristics: Intelligent search strategies for computer problem solving. Addison-Wesley Longman Publishing Co., Boston
35. Guyon I (2002) Gene selection for cancer classification using support vector machines. *Mach Learn* 46(1-3):389. <https://doi.org/10.1023/A:1012487302797>
36. Bach FR (2008) <https://doi.org/10.1145/1390156.1390161>. <http://arxiv.org/abs/0804.1302>
37. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V et al (2011) Scikit-learn: Machine learning in Python. *J Mach Learn Res* 12:2825
38. Sutter JM, Kalivas JH (1993) Comparison of forward selection, backward elimination, and generalized simulated annealing for variable selection. *Microchem J* 47(1-2):60. <https://doi.org/10.1006/mchj.1993.1012>
39. Zou H, Hastie T (2005) Regularization and variable selection via the elastic net. *JR Stat Soc Ser B Stat Methodol* 67(2):301. <https://doi.org/10.1111/j.1467-9868.2005.00503.x>
40. Efron B, Hastie T, Johnstone I, Tibshirani R (2004) Least angle regression. *Ann Stat* 32(2):407. <https://doi.org/10.1214/009053604000000067>. <http://statweb.stanford.edu/tibs/ftp/lars.pdf>
41. Zare H (2015) FeaLect: Scores Features for Feature Selection. <https://cran.r-project.org/package=FeaLect>
42. Gregorutti B, Michel B, Saint-Pierre P (2017) Correlation and variable importance in random forests. *Stat Comput* 27(3):659–678
43. Strobl C, Boulesteix AL, Kneib T, Augustin T, Zeileis A (2008) Conditional variable importance for random forests. *BMC Bioinforma* 9(23):307. <https://doi.org/10.1186/1471-2105-9-307>
44. Toloşi L, Lengauer T (2011) Classification with correlated features: Unreliability of feature ranking and solutions. *Bioinformatics* 27(14):1986. <https://doi.org/10.1093/bioinformatics/btr300>
45. Strobl C, Boulesteix AL, Zeileis A, Hothorn T (2007) Bias in random forest variable importance measures: illustrations, sources and a solution. *BMC Bioinformatics* 8:25. <https://doi.org/10.1186/1471-2105-8-25>. <http://www.ncbi.nlm.nih.gov/pubmed/17254353>