Supplemental Material for

Descriptors for dielectric constants of perovskite-type oxides by materials informatics with first-principles density functional theory

Yusuke Noda^a, Masanari Otake^b and Masanobu Nakayama^{a,b,c,d}

^aCenter for Materials research by Information Integration (CMI²), Research and Services Division of Materials Data and Integrated System (MaDIS), National Institute for Materials Science (NIMS), 1-2-1 Sengen, Tsukuba, Ibaraki 305-0047, Japan

^bFrontier Research Institute for Materials Science (FRIMS), Nagoya Institute of Technology, Gokiso, Showa, Nagoya, Aichi 466-8555, Japan

^cGlobal Research Center for Environment and Energy based on Nanomaterials Science (GREEN), National Institute for Materials Science (NIMS), 1-1 Namiki, Tsukuba, Ibaraki 305-0047, Japan

^dElements Strategy Initiative for Catalysts and Batteries (ESICB), Kyoto University, 1-30 Goryo-Ohara, Nishikyo, Kyoto 615-8245, Japan

Supplemental Section 1. Underlying model for partial least-squares regression

Partial least-squares (PLS) regression is a useful multivariate regression analysis method. Consider linear relationships of *n* samples between explanatory variables x_{ij} (*i* = 1-*n*; *j* = 1-*m*) and objective variables y_i (in this paper, only a single variable y_i is considered):

$$y_i = b_1 x_{i1} + b_2 x_{i2} + b_3 x_{i3} + \dots + b_m x_{im} + f_i$$
(S1),

where b_j and f_i are the coefficients of x_{ij} and the residual, respectively. Eq. S1 can be rewritten using objective variables $\mathbf{y} = (y_1 \ y_2 \ y_3 \ \dots \ y_n)^T$, regression coefficients $\mathbf{b} = (b_1 \ b_2 \ b_3 \ \dots \ b_m)^T$, \mathbf{y} -residuals $\mathbf{f} = (f_1 \ f_2 \ f_3 \ \dots \ f_n)^T$, and explanatory variables $\mathbf{X} = (\mathbf{x}_1 \ \mathbf{x}_2 \ \mathbf{x}_3 \ \dots \ \mathbf{x}_m)$:

$$y = Xb + f \tag{S2}$$

where $x_j = (x_{1j} x_{2j} x_{3j} \dots x_{nj})^T$. In this study, we use the nonlinear iterative partial least squares (NIPALS) [S1] algorithm to predict the regression coefficients between *X* and *y*, which are scaled by a combination of mean centering and standardization (in our

case, the NIPALS algorithm is non-iterative because the objective variables y are not expressed as a matrix but rather a single vector). The objective variables y can be expressed as follows in the NIPALS algorithm:

$$\mathbf{y} = \mathbf{T}\mathbf{c} + \mathbf{f} = \mathbf{X}\mathbf{W}(\mathbf{P}^{\mathrm{T}}\mathbf{W})^{-1}\mathbf{c} + \mathbf{f}$$
(S3),

where $T = (t_1 \ t_2 \ t_3 \ ... \ t_a)$, (*a* indicates the number of factors considered in the PLS regression), $P = (p_1 \ p_2 \ p_3 \ ... \ p_a)$, $W = (w_1 \ w_2 \ w_3 \ ... \ w_a)$, and $c = (c_1 \ c_2 \ c_3 \ ... \ c_a)^T$ are *X*-scores, *X*-loadings, *X*-weights, and *y*-weights, respectively. The first *X*-weight component w_1 and the first *X*-score component t_1 are calculated as follows:

$$w_1 = X^T y / ||X^T y||$$
(S4),
$$t_1 = X w_1$$
(S5).

Component w_a is obtained in order to maximize covariance between X and y. The component t_a is the latent variable used to build the multivariate regression model with a small number of variables. The first X-loading component p_1 and the first y-weight component c_1 can be obtained by a single regression to X and y using the latent variable t_1 . Therefore, the components p_1 and c_1 are also calculated as follows:

$$\boldsymbol{p}_1 = \boldsymbol{X}^{\mathrm{T}} \boldsymbol{t}_1 / \boldsymbol{t}_1^{\mathrm{T}} \boldsymbol{t}_1 \tag{S6},$$

$$c_1 = y^1 t_1 / t_1^{-1} t_1 \tag{S7}.$$

Finally, the projection data to X and y (i.e., $t_1p_1^T$ and t_1c_1) are removed from the original X and y:

$$\boldsymbol{X}' = \boldsymbol{X} - \boldsymbol{t}_1 \boldsymbol{p}_1^{\mathrm{T}}$$
(S8),

$$\mathbf{y}' = \mathbf{y} - \mathbf{t}_1 c_1 \tag{S9}$$

The deflated data X' and y' are used as the next X and y in Eq. S4 to determine the second components t_2 , p_2 , and c_2 . The latent variables t_a are estimated so that they are uncorrelated with each other by repeating the deflation execution to X and y. After calculating the *a*-th principal components of each variable, we can obtain the PLS regression coefficients b^{PLS} with predicted objective variables \hat{y} that satisfy $\hat{y} = Xb^{PLS}$:

$$\boldsymbol{b}^{\text{PLS}} = \boldsymbol{W}(\boldsymbol{P}^{\text{T}}\boldsymbol{W})^{-1}\boldsymbol{c}$$
(S10).

In the PLS method, we use an index to determine the superiority of the regression coefficients. Wold et al. proposed variable importance in projection (VIP) scores,

which evaluate the influence of explanatory variables X on the PLS regression model [S2, S3]. The VIP score for the *j*-th variable is expressed as:

$$VIP_{j} = \sqrt{\frac{\sum_{f}^{F} w_{jf}^{2} \cdot SSY_{f} \cdot J}{SSY_{total} \cdot F}}$$
(S11).

In this equation, w_{if} indicates the weights of the *j*-th variable and *f*-th component. SSY_f is the sum of the squares of *y* explained by the *f*-th component and SSY_{total} is the total sum of the squares of *y* explained by all components; both values are given as follows:

$$SSY_f = (b^{PLS}_f)^2 t_f^T t_f$$
(S12),

$$SSY_{total} = (\boldsymbol{b}^{PLS})^2 \boldsymbol{T}^T \boldsymbol{T}$$
(S13),

where b^{PLS}_{f} is the *f*-th component of the PLS regression coefficients b^{PLS} , *J* is the number of explanatory variables *X*, and *F* is the total number of components. Explanatory variables with larger VIP scores are important for building the PLS regression model.

Supplemental Section 2. Additional PLS regression analysis

Additional PLS regressions were performed to investigate (i) DOS energy alignment and (ii) the inclusion of perovskites that contains Pb^{2+} and Sn^{2+} at A-sites. Due to the lack of bulk modulus data for Pb^{2+} and Sn^{2+} -containing perovskites and data format alignment, we recalculated the entire PLS regression using slightly different datasets of explanatory variables from those shown in Figure 2 in the main text. Therefore, the diagnostic plots shown in Supplemental Figures S4(a) and S9(a), which correspond to the datasets in Figure 2 in main text, shows differences that help to compare the PLS regression results under the same regression conditions.



Supplemental Figure S1. Plot of predicted residual error sum of squares (PRESS) against the number of components in the PLS regression model for test data. The minimum PRESS value occurs at ten components.



Supplemental Figure S2. Diagnostic PLS-regression plots of logarithmic dielectric constant, ln (ε), for the samples only perovskite with vertices-shared BO₆ octahedra, *i.e.* face- and edge-shared perovskites of R3 and P6₃/mmc symmetries are removed from the samples. The resulting statistical evaluation parameters are as follows; RMSE(training) = 0.43, RMSE(test) = 0.57, *R*²(training) = 0.79, and *R*²(test) = 0.60.



Supplemental Figure S3. Diagnostic PLS-regression plots of logarithmic dielectric constant, ln (ϵ), obtained by adding band gap data to the explanatory variables. The resulting statistical evaluation parameters are as follows; RMSE(training) = 0.29, RMSE(test) = 0.54, R^2 (training) = 0.92, and R^2 (test) = 0.84.



Supplemental Figure S4. Diagnostic PLS-regression plots of logarithmic dielectric constant, ln (ε), using the DOS energy scale aligned with: (a) the Fermi level and (b) the O 2s core level for DOS-derived explanatory variables. RMSEs of the training and test data are presented in the plots. (c) VIPs for various descriptors derived from PLS regression. Upper and lower bar graphs correspond to PLS derived VIPs using DOS descriptors whose energies are aligned to the Fermi level and the O 2s core level, respectively.



Supplemental Figure S5. (a) Relationship between DFPT-derived dielectric constants and Shannon's ionic radii of A ions, which has relatively high PLS-VIP scores. (b) Averaged dielectric constant (dot) and corresponding standard deviation (error bar) for each A ion.



Supplemental Figure S6. Relationship between Bader charge of B ions and dielectric constants. (Several compounds with a low dielectric constant ($\varepsilon < 1$) are removed from the figure as outliers.) The line in the figure corresponds to the results of the least-squares linear fitting. The correlation coefficient *R* is -0.45, and the root-mean-square error (RMSE) is 0.40 for the fitted function.



Supplemental Figure S7. Relationship between dielectric constants and interatomic distance between (a) two A sites and (b) two B sites. Panel (c) is a magnification of panel (b) for several specific A and B compositions.



Supplemental Figure S8. Correlation matrix of dielectric constants as functions of two descriptors among the six explanatory variables a–f for all sample data. The dielectric constants (ln ε) are shown by color gradation. Relatively high dielectric materials (red plots) are condensed in some of regions, but the blue plots also coexist in the same region, indicating no significant correlation (one of examples is the correlation graph between B cation DOS and charge difference).



Supplemental Figure S9. Diagnostic PLS-regression plots of logarithmic dielectric constant, ln (ϵ) for the sample sets (a) without and (b) with perovskites that contains Pb²⁺ or Sn²⁺ at A-site. RMSEs of training and test data are presented in the plots. (c) VIPs for various descriptors derived from PLS regression. Upper and lower bar graphs correspond to PLS derived VIPs with and without perovskites that contains Pb²⁺ or Sn²⁺ at A-site in the regression samples, respectively.

References

- [S1] S. Wold, M. Sjöström, and L. Eriksson, PLS-Regression: A Basic Tool of Chemometrics, Chemom. Intell. Lab. Syst. 58, 109-130 (2001).
- [S2] S. Wold, A. Johansson, and M. Cochi, in 3D QSAR in Drug Design, Theory, Methods, and Applications, edited by H. Kubinyi (ESCOM Science Publishers, Leiden, 1993).
- [S3] M. Farrés, S. Platikanov, S. Tsakovski, and R. Tauler, Comparison of the Variable Importance in Projection (VIP) and of the Selectivity Ratio (SR) Methods for Variable Selection and Interpretation, J. Chemom. 29, 528-536 (2015).