# ML Approaches for Predicting Mechanical Properties of Artificial Spider Silk from Spinning Conditions — Supplementary Material

1st Erik Lidbjörk
*Division of Robotics*
*Perception and Learning*
*KTH Royal Institute of Technology*
Stockholm, Sweden
0009-0006-6794-7030, eriklidb@kth.se

2nd Hedvig Kjellström
*Division of Robotics*
*Perception and Learning*
*KTH Royal Institute of Technology*
Stockholm, Sweden
0000-0002-5750-9655, hedvig@kth.se

3rd Neeru Dubey
*Division of Robotics*
*Perception and Learning*
*KTH Royal Institute of Technology*
Stockholm, Sweden
needub@kth.se

TABLE I

HYPERPARAMETER SETTINGS FOR EACH PARAMETER TO THE HISTGRADIENTBOOSTINGREGRESSOR (HGBR) ESTIMATOR IN SKLEARN. WE DISPLAY THE RANGE WHERE THE VALUES WHERE SAMPLED FROM IN THE HYPERPARAMETER TUNING, AS WELL AS THE SAMPLE MEANS AND SAMPLE STANDARD DEVIATIONS BETWEEN EACH HGBR BETWEEN EACH OUTER-FOLD—MECHANICAL PROPERTY PAIR.

| | Range | Values $A$ | Values $B$ |
|---|---|---|---|
| Learning rate | $[0.01, 0.3]$ | 1.02E-01 ± 9.43E-02 | 5.13E-02 ± 3.76E-02 |
| Max leafs per tree | $[10, 50]$ | 1.71E+01 ± 9.54E+00 | 1.86E+01 ± 1.08E+01 |
| Minimum samples per leaf | $[1, 20]$ | 9.73E+00 ± 7.82E+00 | 8.93E+00 ± 6.88E+00 |
| $L_2$ regularization | $[0.001, 10]$ | 2.23E+00 ± 3.10E+00 | 1.79E+00 ± 2.64E+00 |
| Maximum iterations | $[0, \infty)^a$ | 9.11E+01 ± 1.49E+01 | 9.62E+01 ± 1.31E+01 |

[a]We used early_stopping=True when training the inner-fold estimators, and took the resulting mean of max_iters for the respective outer-fold estimator.
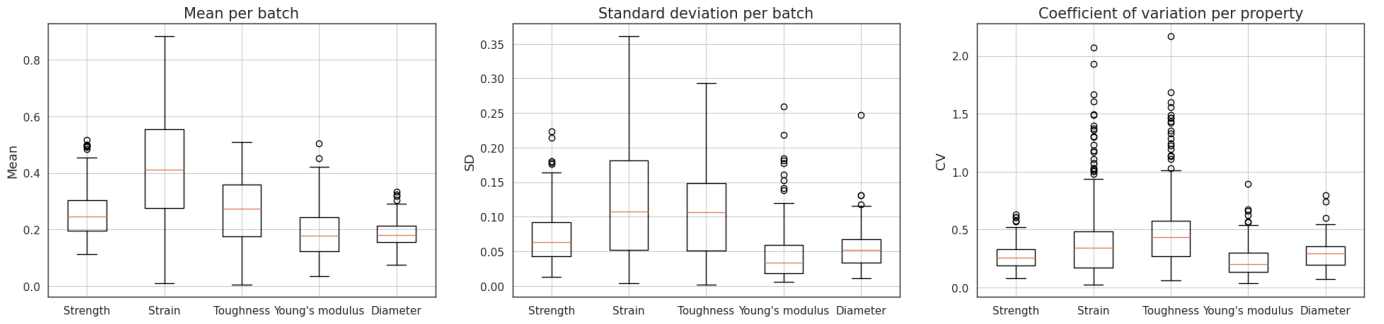


Fig. 1. For each spinning condition, around ten measurements of the mechanical properties were made. This figure depicts the mean, sample standard deviation (SD), and coefficient of variation (mean/std) (CV) of the mechanical properties, per spinning condition—targets batch, as a way to depict label noise in the mechanical property measurements in batches. Each property has batches with CV¿50%, and some over 200% (for strain and toughness), which indicates high variability compared to means.
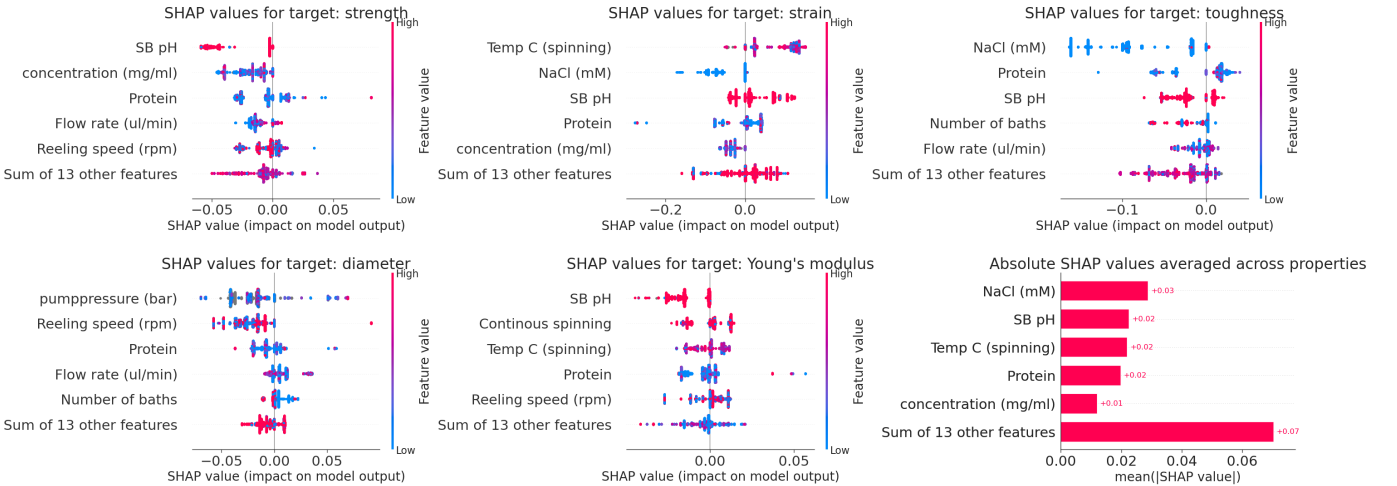


Fig. 2. Top five feature importance ranking for each mechanical property, derived from model $A$. Ranking is decided based on the total sum of absolute value SHAP-values for each feature. A high/low SHAP-value (position on the $x$-axis) indicates positive/negative correlation with model predictions. Colors indicate the numerical value of the feature (red=high, blue=low).
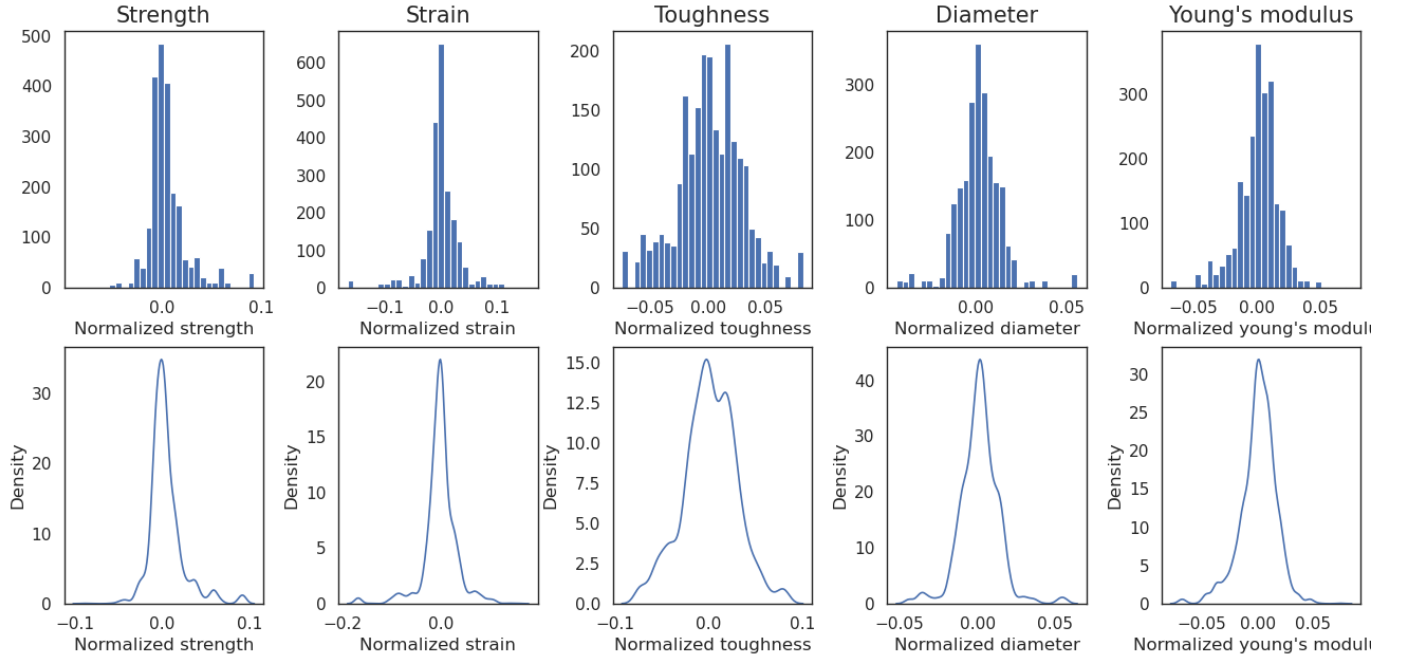
Fig. 3. The distribution between differences in mechanical property predictions between models $A$ and $B$, depicted as a histogram and kernel density estimation (KDE). We deem these distributions to be approximately normal, due to their symmetry, low skew, and bell curve shape. We therefore compare the models using the Student's paired $t$-test.
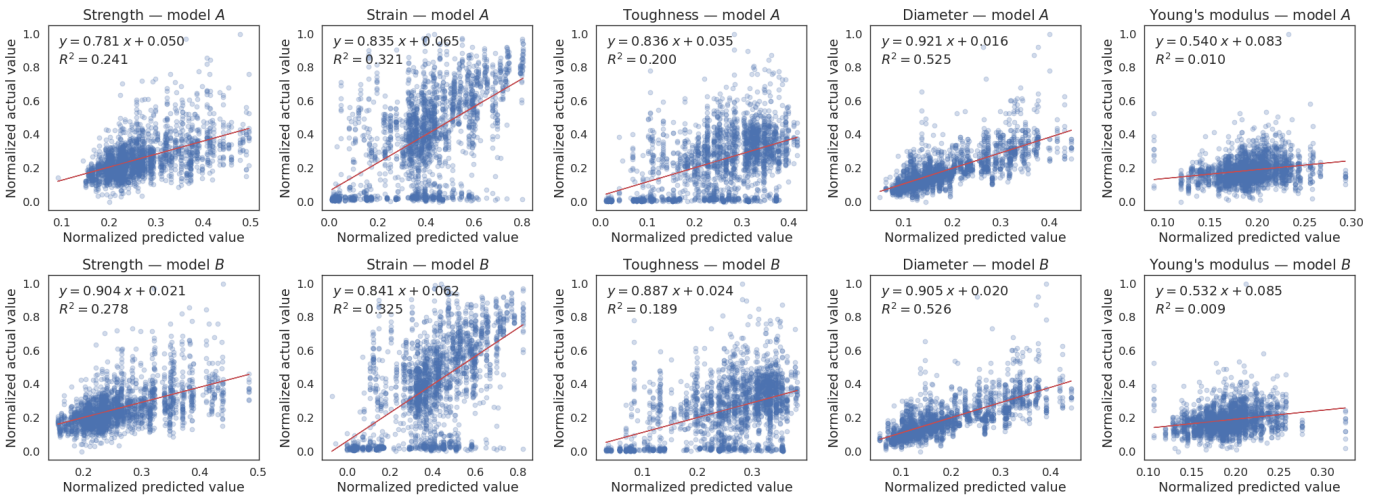


Fig. 4. Ground-truth ($y$-axis) vs. predicted ($x$-axis) values for each mechanical property, depicted for both models $A$ and $B$. For the following graphs, we have a slope coefficient $0 < k < 1$. We clearly see outliers present in the dataset ignored by our models, evident by either a high or low ground truth value on the $y$-axis-axis. This is most apparent for strain and toughness, where we have an abundance of points around zero on the $y$-axis, while the models predict higher.
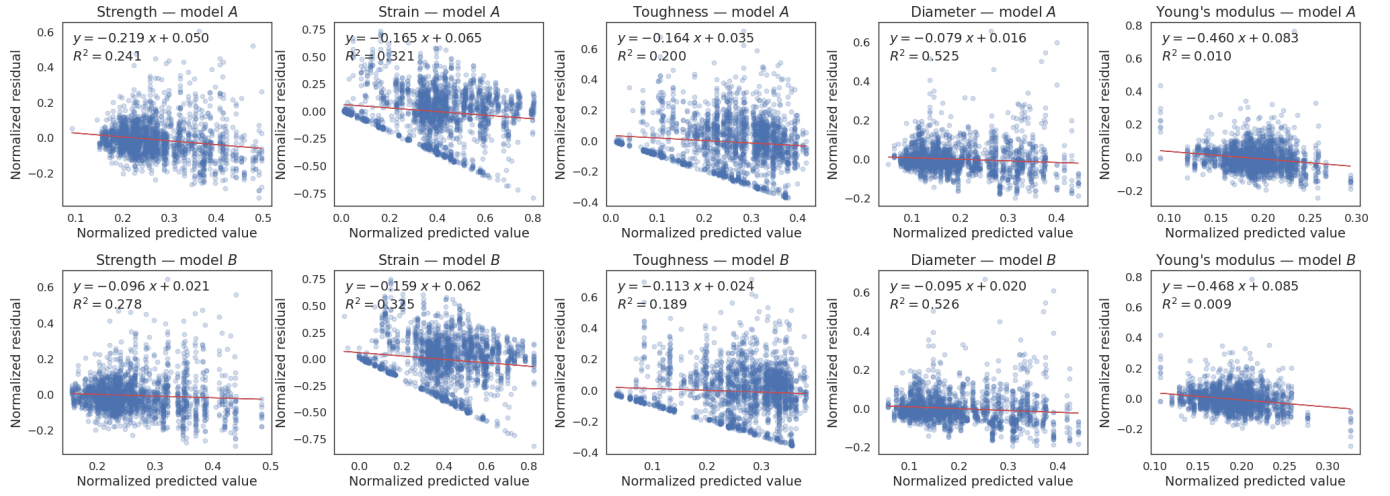
Fig. 5. Residuals (ground-truth - predicted value) ($y$-axis) vs. predicted ($x$-axis) values for each mechanical property, depicted for both models $A$ and $B$. For the following graphs, we have a slope coefficient $k < 0$. The phenomenon where an abundance of zero values is predicted higher by the models for strain and toughness, shown in Fig 4, can here be seen by a line formed by linearly increasing residuals.