# Ancient Artifacts

Team: Weixi Chen, Shea Conaway, Shuyang Lin, Sydney Simmons

Team Lead: Mubarak Ganiyu

# Agenda

**Project Objective:**

Identify the locations of Mayan stone tool manufacturing areas by characterizing the composition of soil samples from ancient villages

- Ancient stonemakers would remove large debris from sites, but small particles would remain, called lithic microdebitage

- Our team trained a model to distinguish between what particles are stone and soil

- Using our model, future particle samples can be analyzed to identify stone debris and predict manufacturing areas

- Identifying tool manufacturing sites can give valuable insight into past cultural activity in Mayan villages
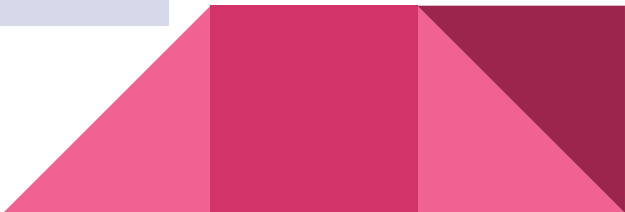
**Initial Datasets**

- Initial data included two datasets (one for stone and one for soil samples) in order to train our model to recognize the differences between the particles and **identify stone**
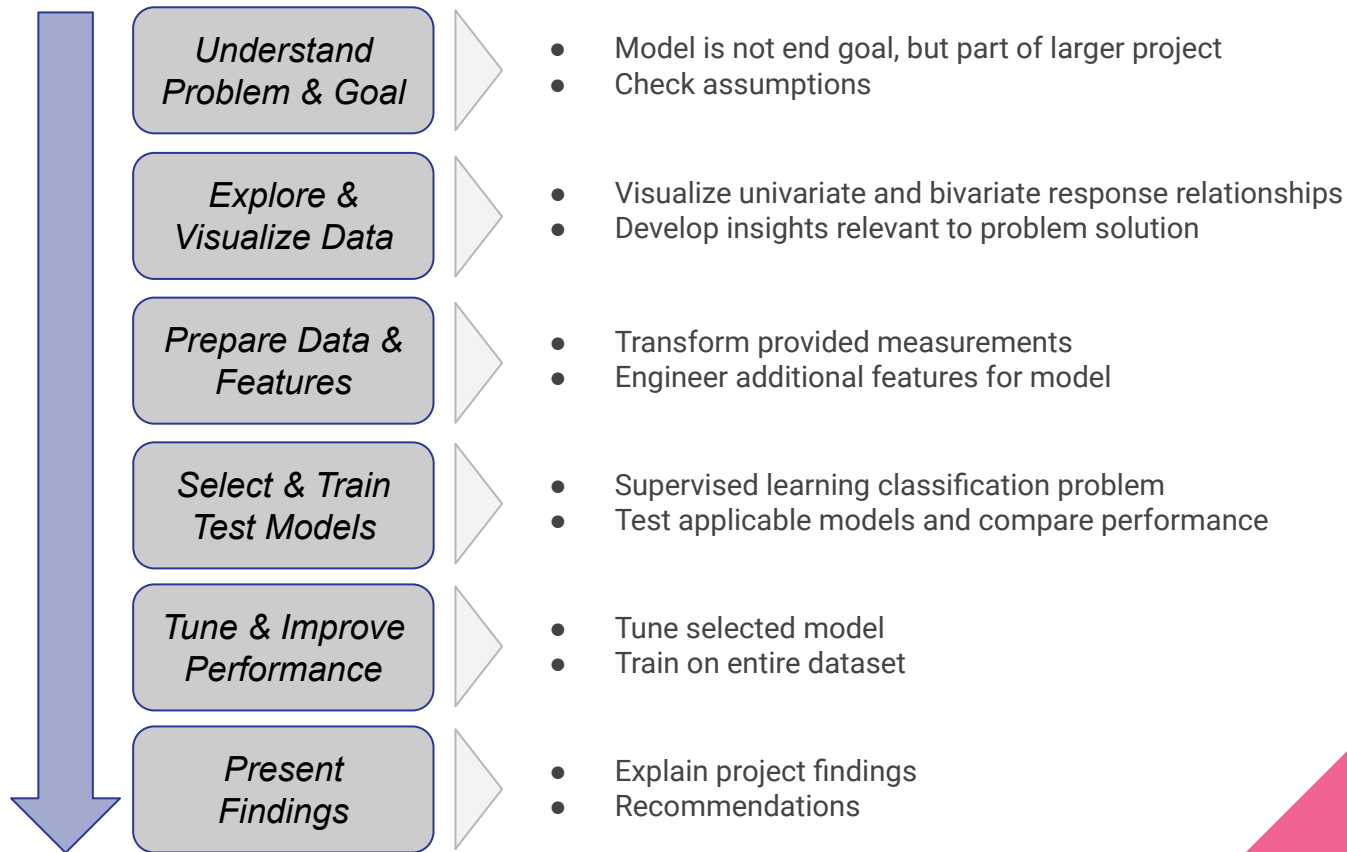
**Data Attributes**

- Combined dataset of both stone and soil included over **75,000 rows of data and 48 features**
- Features were numeric measurements and ratios to describe the particles, including transparency, volume, area, and thickness

**Splitting Data**

- After aggregation and data cleaning, the team split the data into data to use to train the model and data to use to test the performance of the model to ensure an **appropriate product for future use**
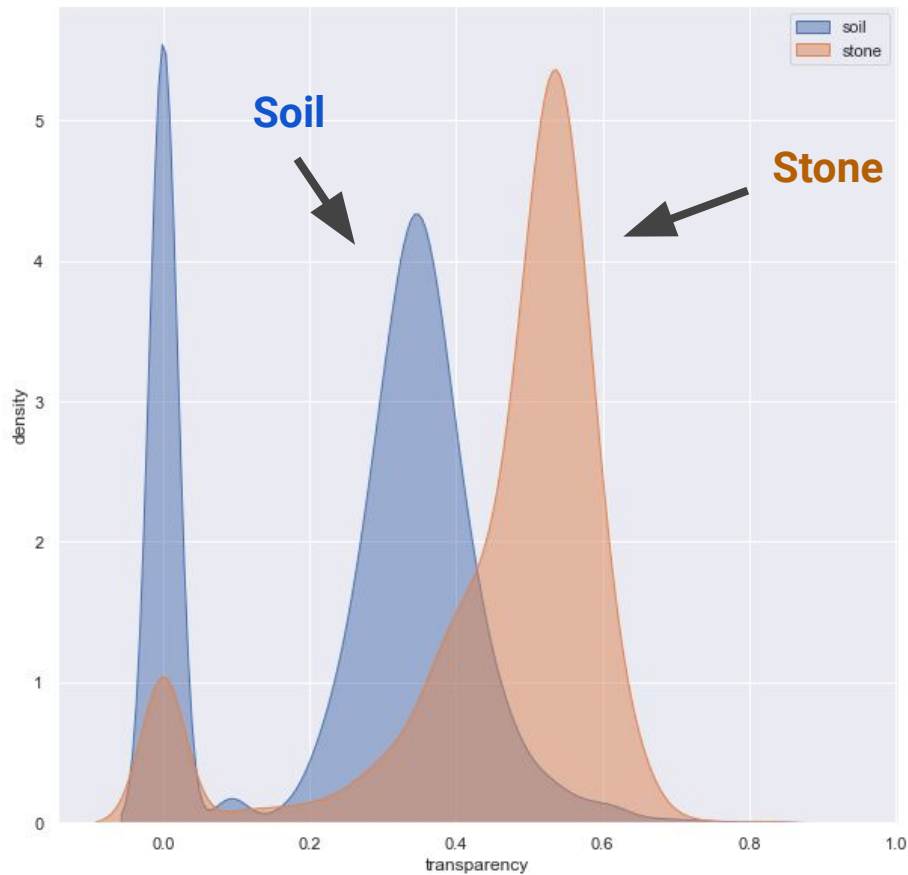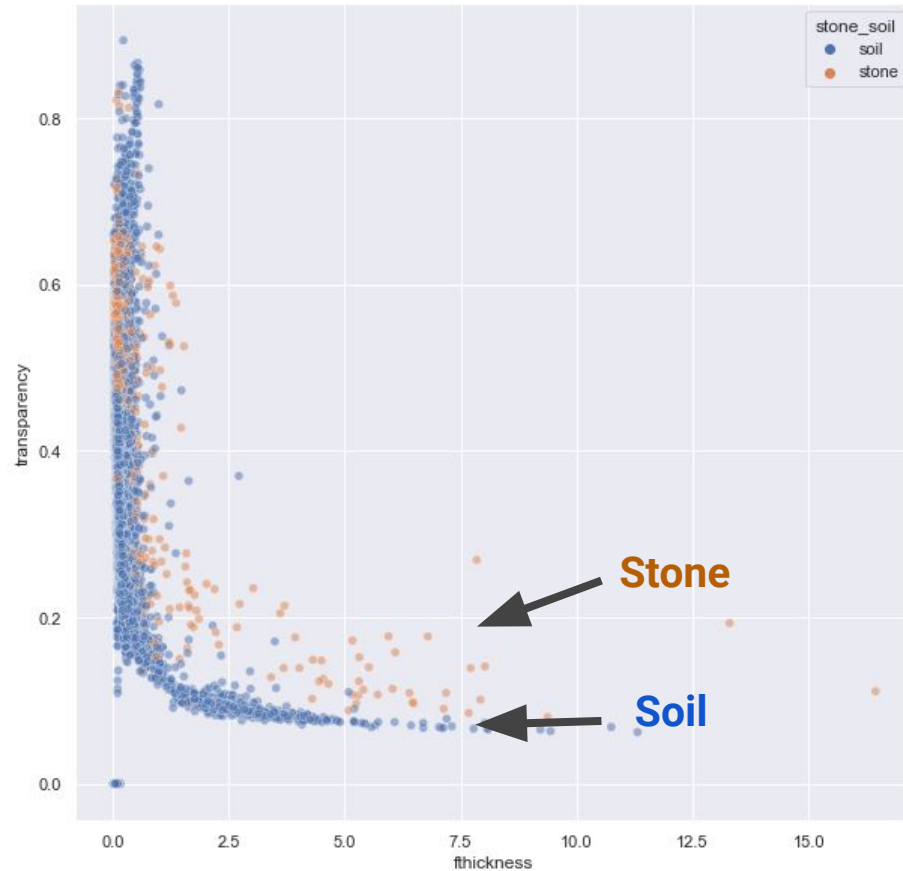
# Data Science Approach

**3**

**Understand Problem & Goal**
- Model is not end goal, but part of larger project
- Check assumptions

**Explore & Visualize Data**
- Visualize univariate and bivariate response relationships
- Develop insights relevant to problem solution

**Prepare Data & Features**
- Transform provided measurements
- Engineer additional features for model

**Select & Train Test Models**
- Supervised learning classification problem
- Test applicable models and compare performance

**Tune & Improve Performance**
- Tune selected model
- Train on entire dataset

**Present Findings**
- Explain project findings
- Recommendations

# Explore, Visualize, and Engineer Features



Transparency

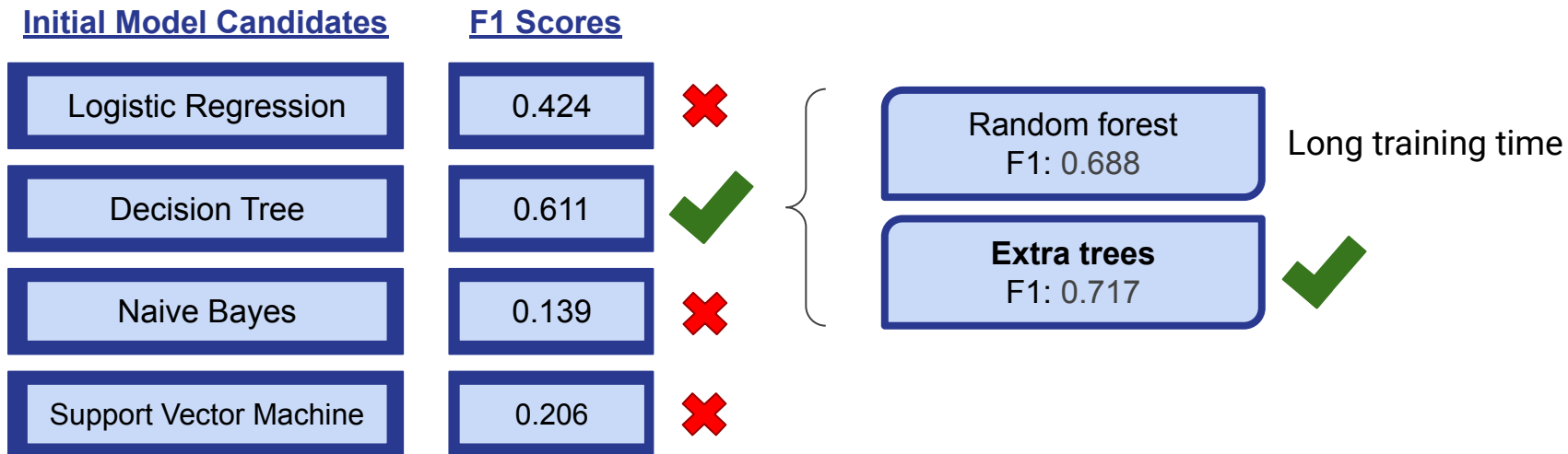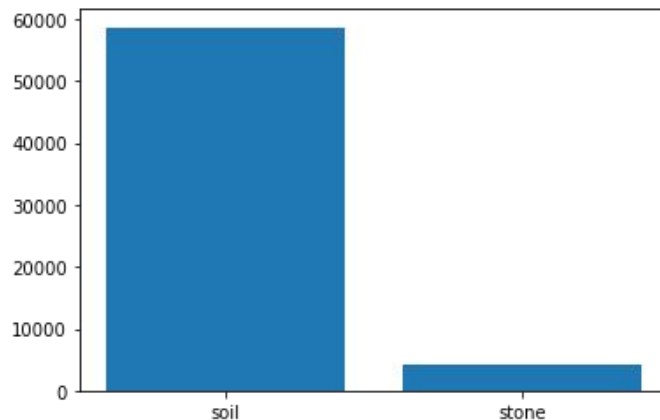Transparency x Thickness

# Key Findings: Model Improvement (1/2)

- We selected **Extra Trees Classifier** as our final model model to limit overfitting and maintain a fast, effective model for the end user.

- In order to improve our model, we also:

  - **Added more stone particles to our training data**

  - Limited the used features using variable importance

- Reorganized the train and test size with respect to stone and soil.

- Used grid search cross validation to look for better hyper-parameters.

- Increased random state and applied bootstrap to improve out-of-bag performance.
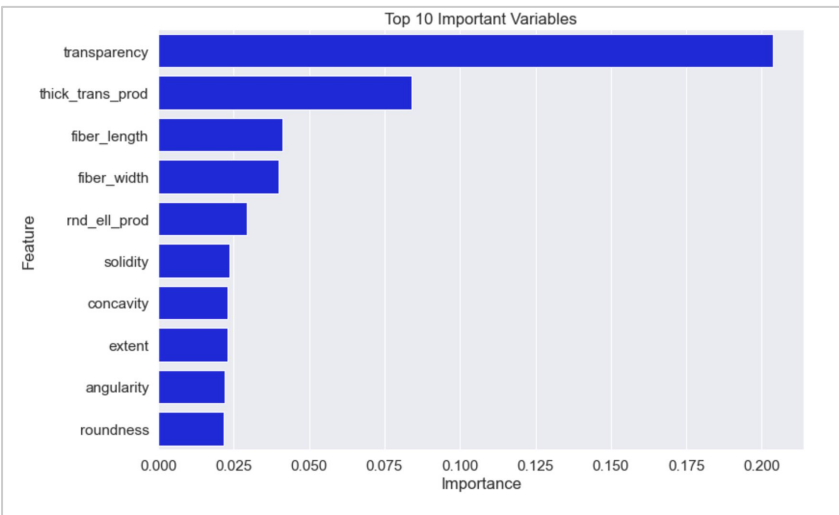
### Histogram of Stone or Soil

# Key Findings: Model Improvement (2/2)

● We selected **Extra Trees Classifier** as our final model model to limit overfitting and maintain a fast, effective model for the end user.

● In order to improve our model, we also:

  ○ Added more stone particles to our training data

  ○ **Limited the used features using variable importance**

● Exported the variable importance table from the initial model

● Figured out the importance of size and transparency

● Attempted to do a sub-model on samples with a larger size
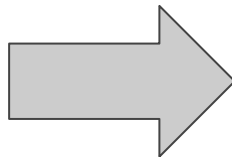


Top 10 Important Variables

# Key Findings: Model Performance

- In improving our model, we focused on increasing **recall and F1 score**, limiting false positives and negatives and increasing true positives (when the model predicted stone and the particle was stone)

| Original Extra Trees Model |
|---|

|  | 14525 | 112 |
|---|---|---|
|  | TN | FP |
|  | FN | TP |
|  | 498 | 587 |
|  | 0 | 1 |

Accuracy: **0.96**
Precision: **0.84**
Recall: **0.54**
F1: **0.66**

| Improved Extra Trees Model |
|---|

|  | 14564 | 99 |
|---|---|---|
|  | TN | FP |
|  | FN | TP |
|  | 671 | 919 |
|  | 0 | 1 |

Accuracy: **0.95**
Precision: **0.91**
Recall: **0.59**
F1: **0.72**

# Next Steps and Recommendations

④

- Sample sizes need to be large enough to have <u>different particle sizes</u>

- Soil / stone determination most accurate on <u>big particles with large perimeters</u>

- Constant <u>model refinement</u> and use to locate stone sites

  - Produce more samples for training

  - Evaluate real-world confusion manually for feedback and improve model

  - Monitor particle analyzer performance and continue to decrease false negative results