

5 cm

Google NACME AMLI Bootcamp

Classifying Microdebitage with Machine Learning

Presenters: Jose Cruz, Kimi Medina-Castellano,

Luke Taylor, Rodrigo Aguilar Barrios

Under mentorship of Dr. Phyllis Johnson and Dr. Reyes-Centeno



The Team



Rodrigo Aguilar Barrios

UC Berkeley
B.S. Electrical
Engineering and
Computer Science



Luke Taylor

University of Kentucky
B.S Computer Science
Digital Media and Design



Kimi Medina-Castellano

University of Kentucky
B.S. Mathematics
Physics & CS



Jose Cruz

University of Kentucky
B.S. Computer Science &
Mathematics

Overview

01.
Background

02.
Data Analysis

03.
Modeling

04.
Future Steps

05.
Acknowledgements

What is Microdebitage?



Leftover pieces (usually 6 mm or smaller) stones
pieces that are knocked off from making stone tool

Tools have separate stages of production

Requires precision and patience

Knapping

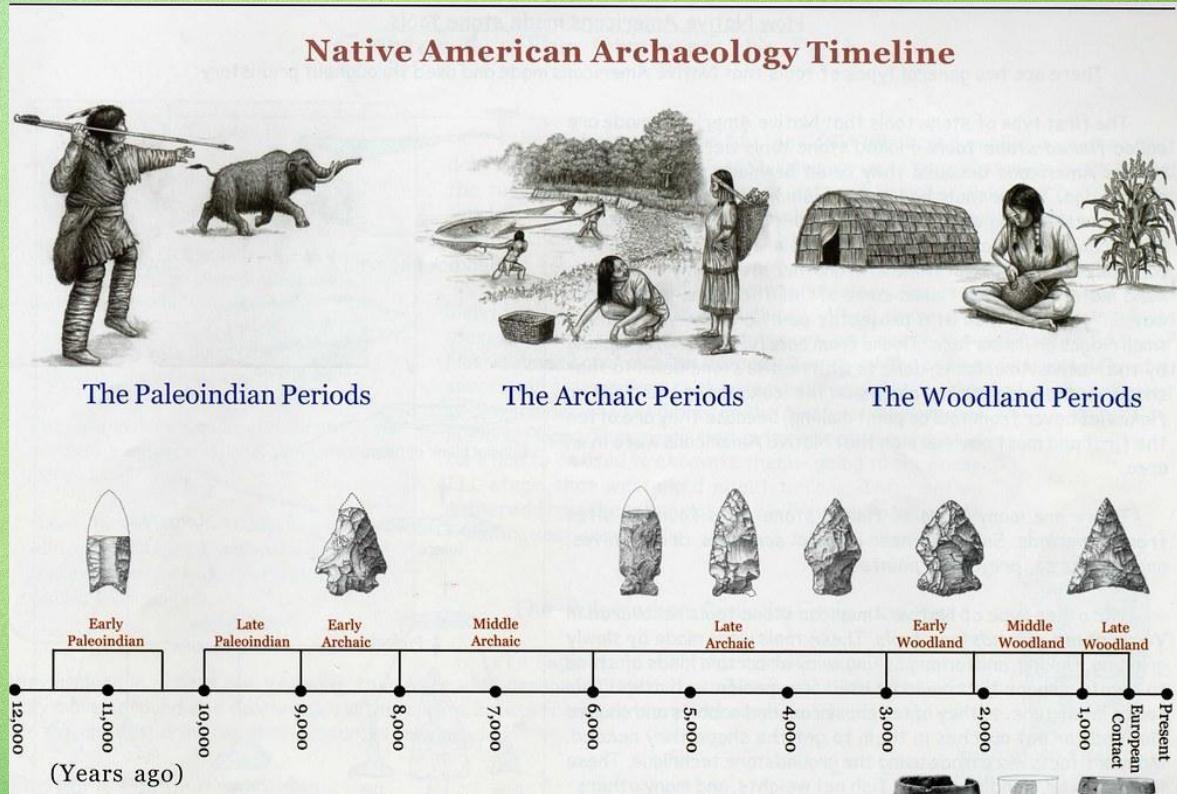


Tools can be recovered due to stone's resilience
from human activity and natural processes



Chert and Obsidian







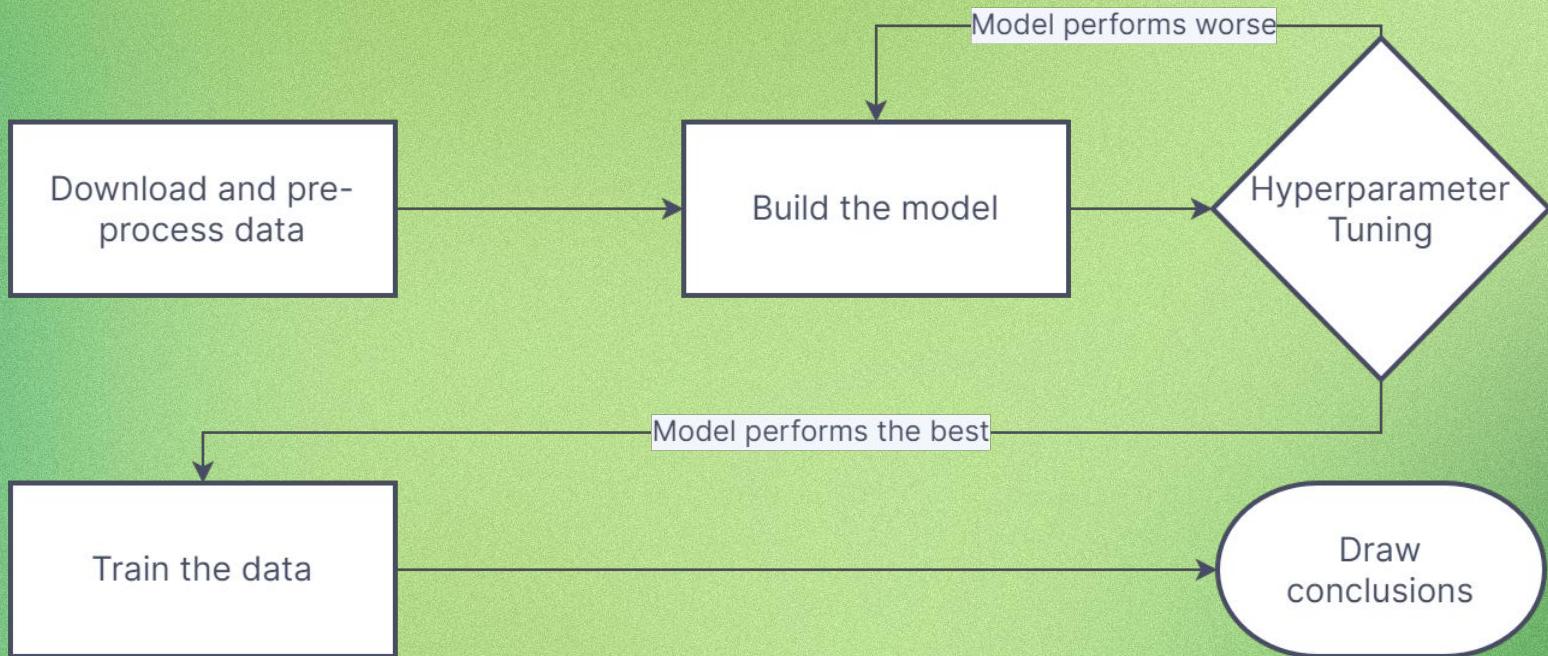
Problem to solve

Takes **>10 hours per sample** for researchers to visually analyze and sort the microdebitage

Solution: Create a model to classify microdebitage into certain stages of tool production

Why? To learn more about ancient civilizations

Procedure



A professional knapper at Vanderbilt University makes simulants of stonedebitage



Knap

Microdebitage is gathered between stages



Gather Microdebitage

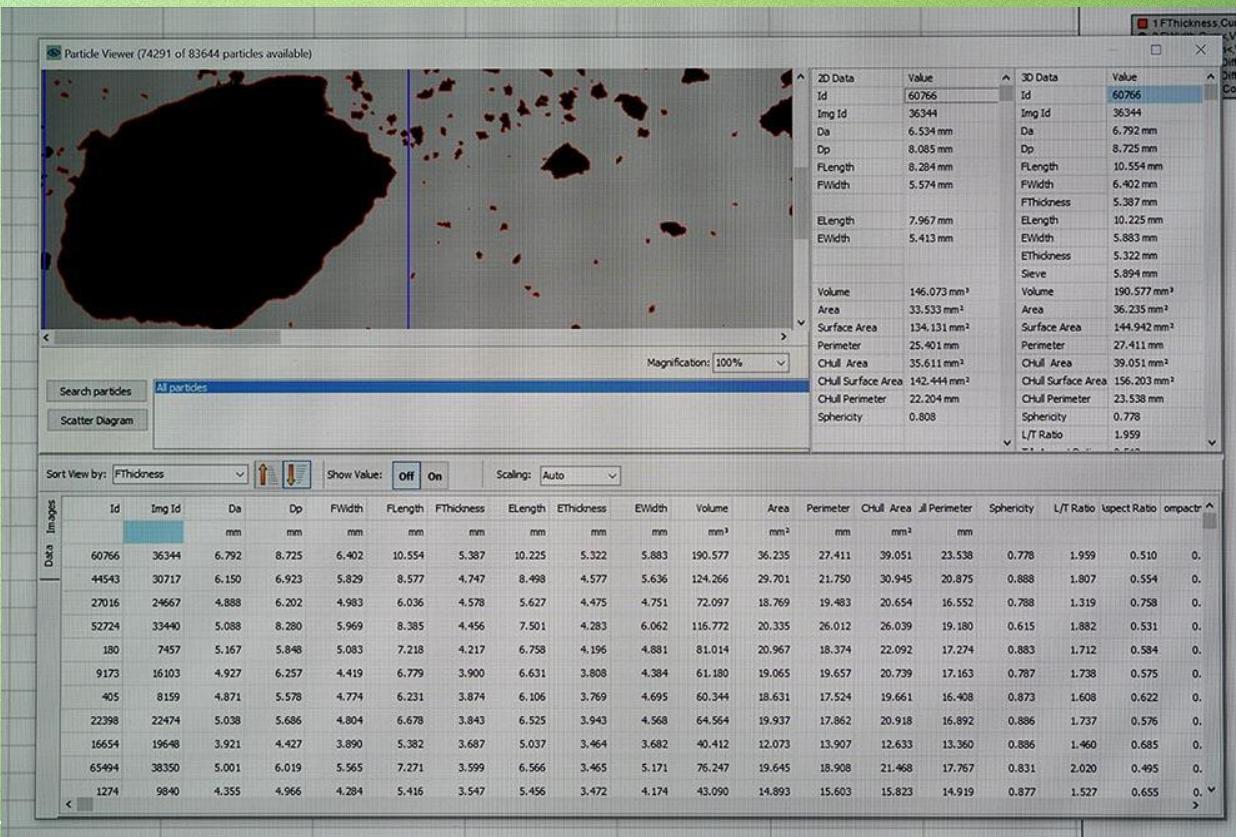
Acquiring Data

Sort through machine

Simulants are ran through Microtrac machine that takes images and numerical measurements of each particle



Data Analysis



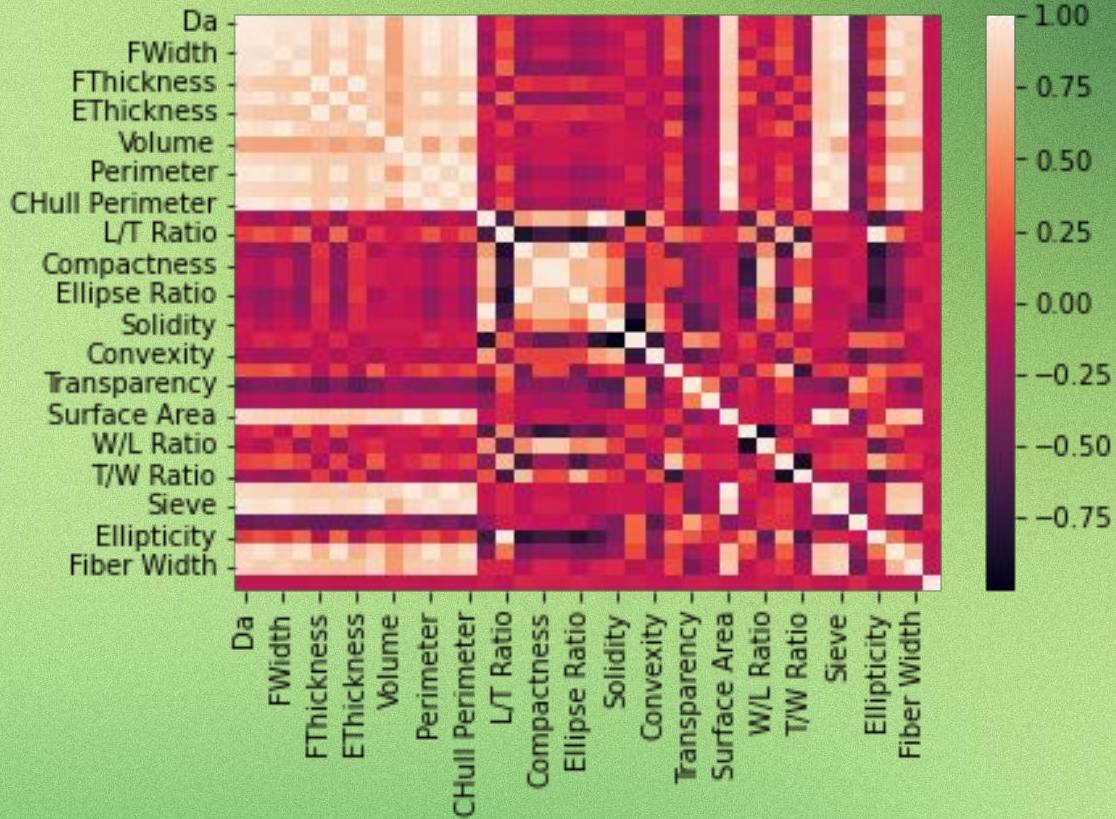
Preprocessing

The features that were taken out were the ones that had repeated values, such as the Filters and the hash columns. The Img Id and Id columns were deemed as not relevant to the model.

There were three columns that one of the data frames had, so we removed those columns from that particular dataframe

	Id	Filter0	Filter1	Filter2	Filter3	Filter4	Filter5	Filter6	hash	Img Id	Curvature	Transparency	Angularity
0	5903	Reject	0	14271	0	0.916	88.889						
1	3002	Reject	0	11576	0.467	0.835	87.5						
2	9893	Reject	0	18211	0.91	0.807	81.667						
3	3713	Reject	0	12265	2.003	0.787	72.381						
4	7862	Reject	0	16105	0	0.786	82.5						
...
36676	7	Reject	0	4022	0	0.154	18.989						
36677	18	Reject	0	5347	0	0.154	21.2						
36678	33	Reject	0	7645	0	0.153	18.144						
36679	38	Reject	0	8347	0	0.148	17.347						
36680	29	Reject	0	7032	0	0.144	29.632						

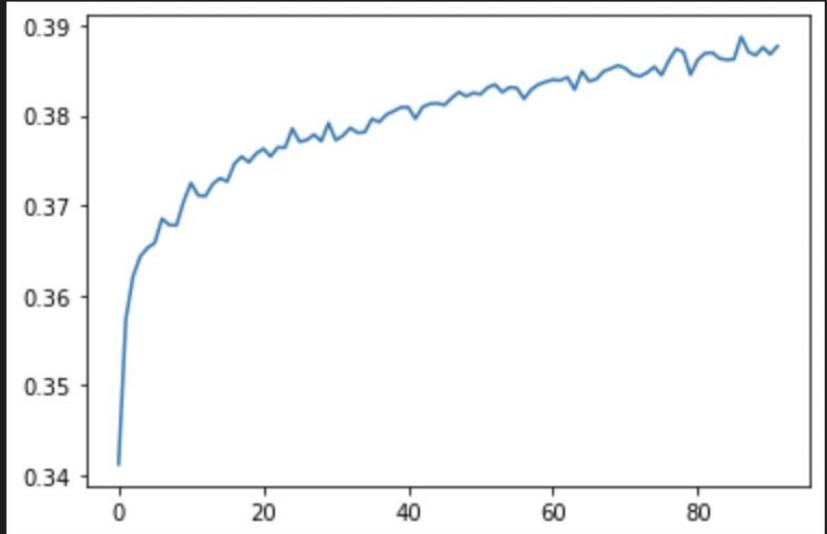
Data Visualization





Tensorflow Multiclass Classification

```
tf.keras.Sequential([  
    tf.keras.layers.Dense(len(filtered) * 2,  
input_shape=(len(filtered),)),  
    tf.keras.layers.Dense(16, activation =  
'relu'),  
    tf.keras.layers.Dense(8, activation =  
'relu'),  
    tf.keras.layers.Dense(5, activation =  
tf.nn.softmax)])  
  
model.compile(  
    loss='categorical_crossentropy',  
    optimizer='Adam',  
    metrics=['accuracy'],  
)
```



Graph of accuracy per epoch
loss: 1.4093 - accuracy: 0.3877
val_loss: 1.4085 - val_accuracy: 0.3884

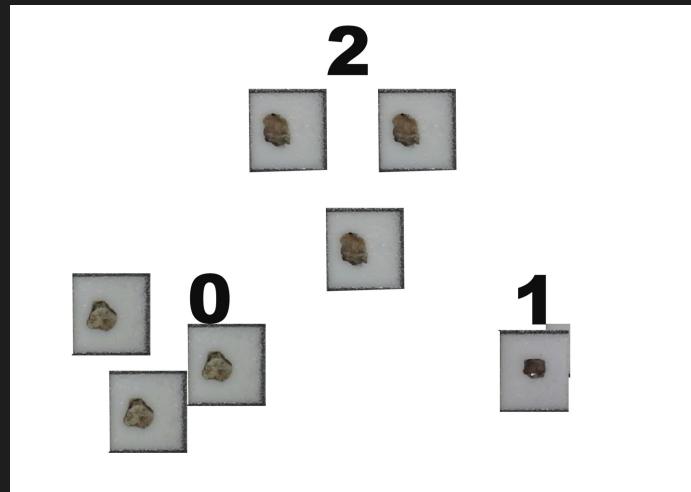
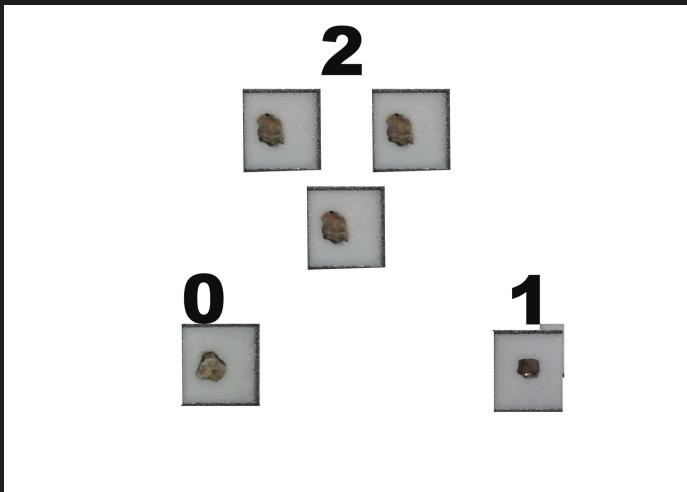


Modeling

KNeighborsClassifier

```
KNeighborsClassifier(n_neighbors=40, weights= 'distance', algorithm= 'brute', leaf_size = 35)
```

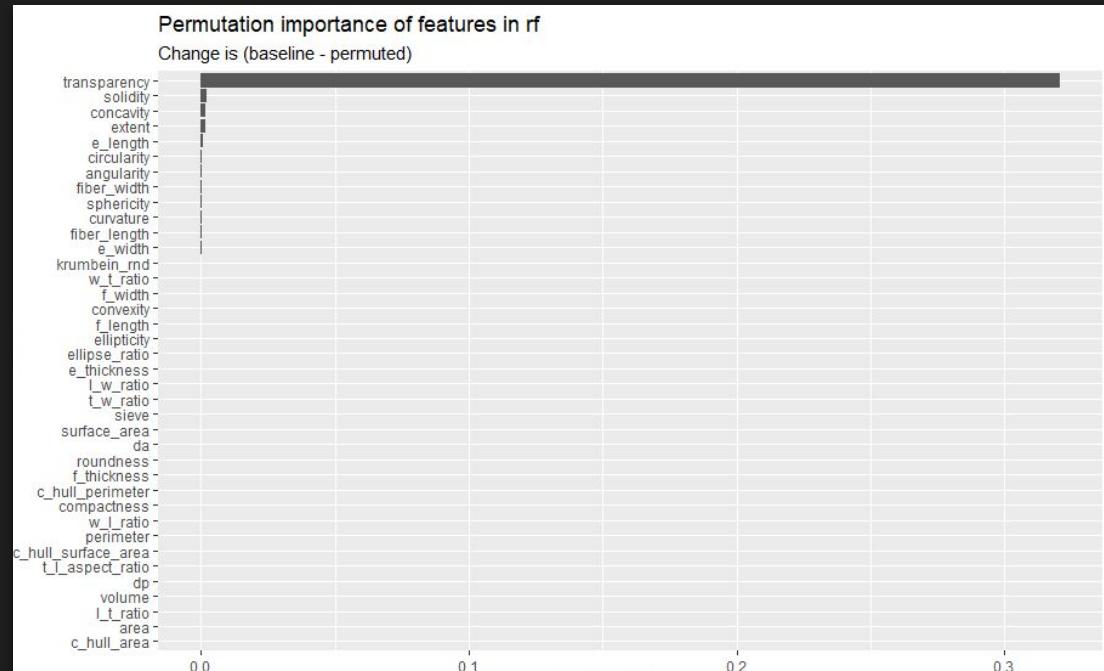
Accuracy: 0.3406



Importance of weights = distance

Making two models

	Angularity	Transparency	Curvature
0	88.889	0.916	0.0
1	87.5	0.835	0.467
2	81.667	0.807	0.91
3	72.381	0.787	2.003
4	82.5	0.786	0.0
...
49489	17.639	0.195	0
49490	25.781	0.306	0.184
49491	39.063	0.232	0
49492	48.889	0.389	0.37
49493	48	0.39	0.168



Two Model K-Neighbor Classifier

Modeling



```
O_model = KNeighborsClassifier(n_neighbors=5,
                               weights= 'distance',
                               p = 1,
                               Leaf_size= 5,
                               algorithm='auto')

O_model.fit(O_data[O_filtered], O_data["Production Stage"])
labels = O_data['Production Stage']

random = O_data.sample(frac = 1)

one_hundred = random[:10000]

O_predictions = O_model.predict(X_test)

print('Accuracy score:',accuracy_score(y_test, O_predictions))
```

```
C_model = KNeighborsClassifier(n_neighbors=5,
                               weights= 'distance',
                               p = 1, leaf_size= 5,
                               algorithm='auto')

C_model.fit(C_data[C_filtered], C_data["Production Stage"])
labels = C_data['Production Stage']

random = C_data.sample(frac = 1)

one_hundred = random[:10000]

C_predictions = C_model.predict(CX_test)

print('Accuracy score:',accuracy_score(cy_test, C_predictions))
```

Accuracy score: 0.7373737373737373

Accuracy score: 0.4396887159533074

Decision Tree for Obsidian Model

Modeling



```
#Creating a list of the features and target
FEATURES = df.columns[:-1]

TARGET = df.columns[-1]

#Splitting our data for training and testing
x_train, x_test, y_train, y_test = train_test_split(
    df[FEATURES],
    df[TARGET],
    test_size = 0.2,
    random_state = 10
)
```

Accuracy: 0.8505909687847257

```
#Creating a decision tree object
dt = tree.DecisionTreeClassifier(
    criterion = 'entropy',
    max_depth = 5,
    random_state = 10
)

#Training our data
dt.fit(x_train, y_train)

#Making Predictions
prediction = dt.predict(x_test)

#Calculating our accuracy
acc = metrics.accuracy_score(y_test, prediction)
print('Accuracy: ', acc)
```



01.

```
# features and target
features = chert_df.loc[:, chert_df.columns != 'Stage']
feature_names = features.columns
```

02.

```
# split data
x_train, x_test, y_train, y_test = train_test_split(
    features,
    target,
    test_size = 0.3,
    random_state = 10
)
# train model
```

03.

```
decision_tree = DecisionTreeClassifier()
decision_tree.fit(x_train, y_train)
```

```
dt = tree.DecisionTreeClassifier(
    criterion = 'gini'
)
dt.fit(x_train, y_train)
y_pred = dt.predict(x_test)
acc = metrics.accuracy_score(y_test, y_pred)
print('Accuracy: ', acc)
```

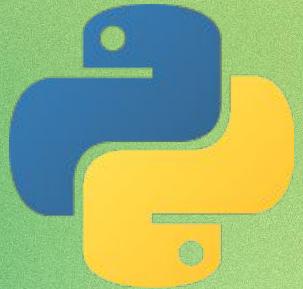
Decision Tree: Chert

The chert model contained 3 datasets with 3 unique stages.

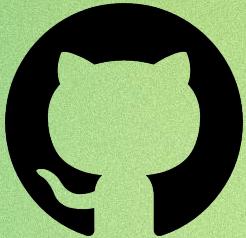
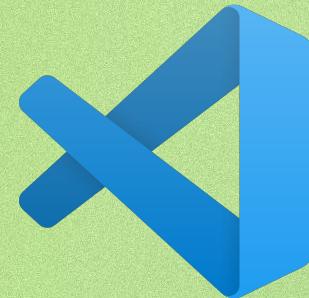
Accuracy: 43.75%

Will be tuning this model to increase accuracy

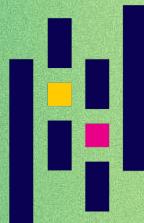
Medium accuracy can be due to missing columns



colab



TensorFlow

 pandas

The pandas logo, featuring a dark blue 'pandas' wordmark next to a stylized icon composed of vertical bars in dark blue, navy, yellow, and pink.

 scikit
learn

The scikit-learn logo, which includes a blue circle and an orange rounded rectangle containing the 'scikit' and 'learn' text.

Obstacles

We were missing three crucial columns for some of our training datasets forcing us to get rid of those columns across all the data.

Running GridSearch to tune models on all the models took a very long time

Angularity	Transparency	Curvature
88.889	0.916	0
87.5	0.835	0.467
81.667	0.807	0.91
72.381	0.787	2.003
...
18.989	0.154	0
21.2	0.154	0
18.144	0.153	0
17.347	0.148	0
29.632	0.144	0

Ethical Consideration

- Misclassifying the time period
- Incorrect assumptions about time period
- Only looking at two materials
- Obtaining microdebitage

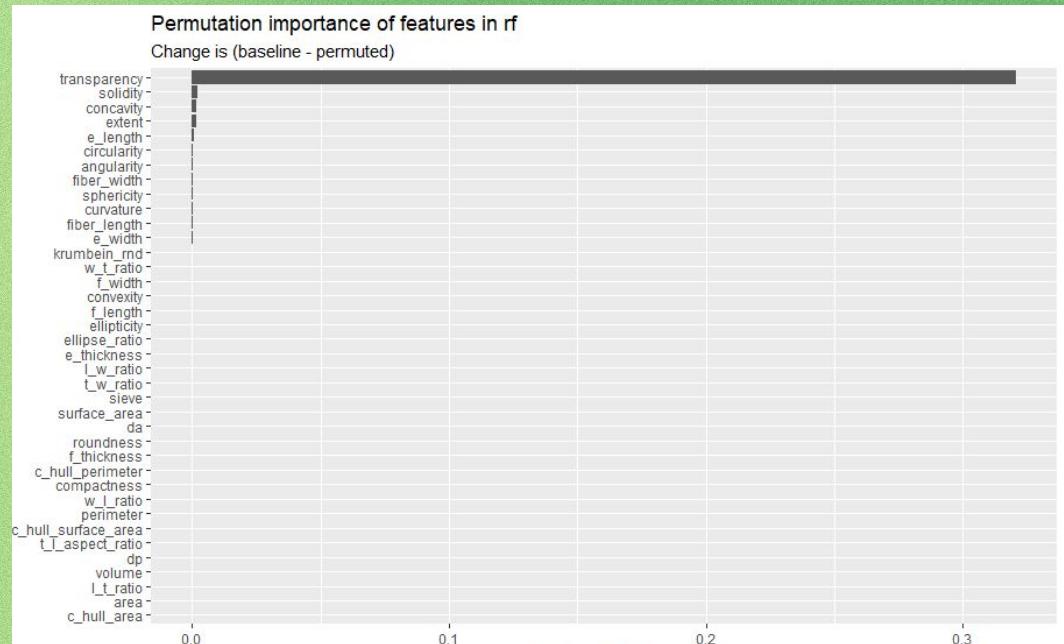
Moving Forward

Combining the chert model with the obsidian model

Combining the model that separates the soil and microdebitage with ours.

Re-analyzing the datasets that were missing columns to fill those data points and add those features into our training for an even better result.

Continuing research with mentors.



Acknowledgements

Dr. Corey Baker

Dr. Phyllis Johnson and mentors

University of Kentucky

NACME

Oscar Feliz

Jalen Collins

References

- <https://towardsdatascience.com/how-to-tune-a-decision-tree-f03721801680>
- <https://scikit-learn.org/stable/>
- <https://www.xoriant.com/blog/decision-trees-for-classification-a-machine-learning-algorithm>
- <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>
- https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
- Eberl, M., Johnson, P. , & Estrada Aguila, R. (2022). Studying lithic microdebitage with a dynamic image particle analyzer. *North American Archaeologist*, 1-16. DOI: 10.1177/01976931221109301
- Unpublished articles provided by Dr. Phyllis Johnson



Questions?

