# "Hominid Chronology and Classification: Machine Learning Applications in Paleoanthropology"

• • •

# Challenges in Classification and Dating in Paleoanthropology

One of the biggest challenges in paleoanthropology is identifying and classifying discovered fossils. Hominids exhibit significant variation within the same species, which can lead to confusion about whether certain fossils represent a new species or simply intraspecific diversity.

Another major challenge is fossil dating. Traditional techniques, such as radiocarbon dating or isotope analysis of surrounding rocks, are often expensive, invasive, or limited by the availability of suitable materials. In this context, data science techniques that correlate physical traits with known time periods can provide a non-invasive and cost-effective tool to complement these traditional methods.

# The Project

This project aims to develop predictive models that classify hominid fossils in terms of their Genus and Species and estimate their geological age using variables related to their physical characteristics. This serves as a preliminary analytical tool to complement traditional paleoanthropological research techniques.

Objectives:

- Taxonomic Classification of Fossils:
  Design and implement supervised classification models that utilize physical traits—such as cranial size, bone length, and endocranial capacity—to accurately predict the Genus and Species of fossils.
- Chronological Estimation:
  Develop regression models to estimate the geological age (approximate time period) of hominid fossils based on the same physical characteristics or derived variables.

# The Data

For this project, we used a hominid dataset sourced from Kaggle, which provides a comprehensive collection of data related to the physical, contextual, and chronological characteristics of fossil specimens. This dataset was chosen due to its richness in relevant information for the project's objectives, including taxonomic classification (Genus_Species) and chronological estimation (Time).

The dataset contains 12,000 records and 27 columns, offering a robust framework for analysis and modeling. Each row represents a fossil specimen, while the columns reflect various associated attributes.

Dataset used in this study:
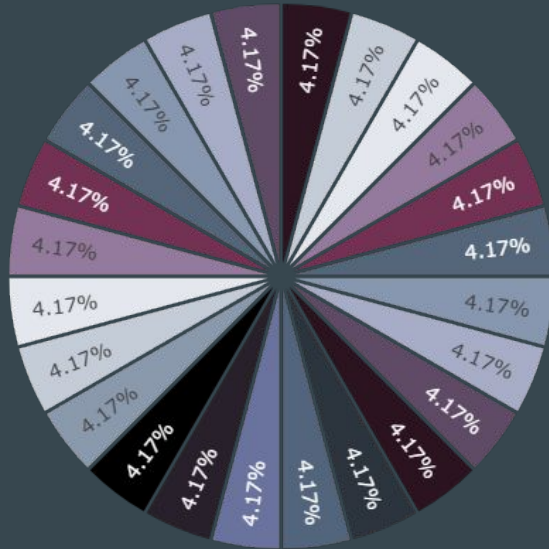🔗 Kaggle - Biological Data of Human Ancestors

# Data Cleaning

The data cleaning process for this project did not require imputing missing values or handling null entries, as the dataset contained no incomplete records. However, several key transformations were performed to enhance data consistency and usability:

- Column Name Standardization:
  Column names were modified to ensure consistency and clarity, using a uniform and descriptive naming convention. This improved the readability and handling of variables during analysis.
- Time Variable Transformation:
  The time variable was standardized to always be expressed as a numerical value representing the number of millions of years before the present.
- Previously, time was recorded in various formats, including ranges and textual descriptions, which complicated the analysis.
- Cranial Capacity Standardization:
  The cranial capacity variable was normalized to ensure all values were expressed in cubic centimeters ($cm^3$).
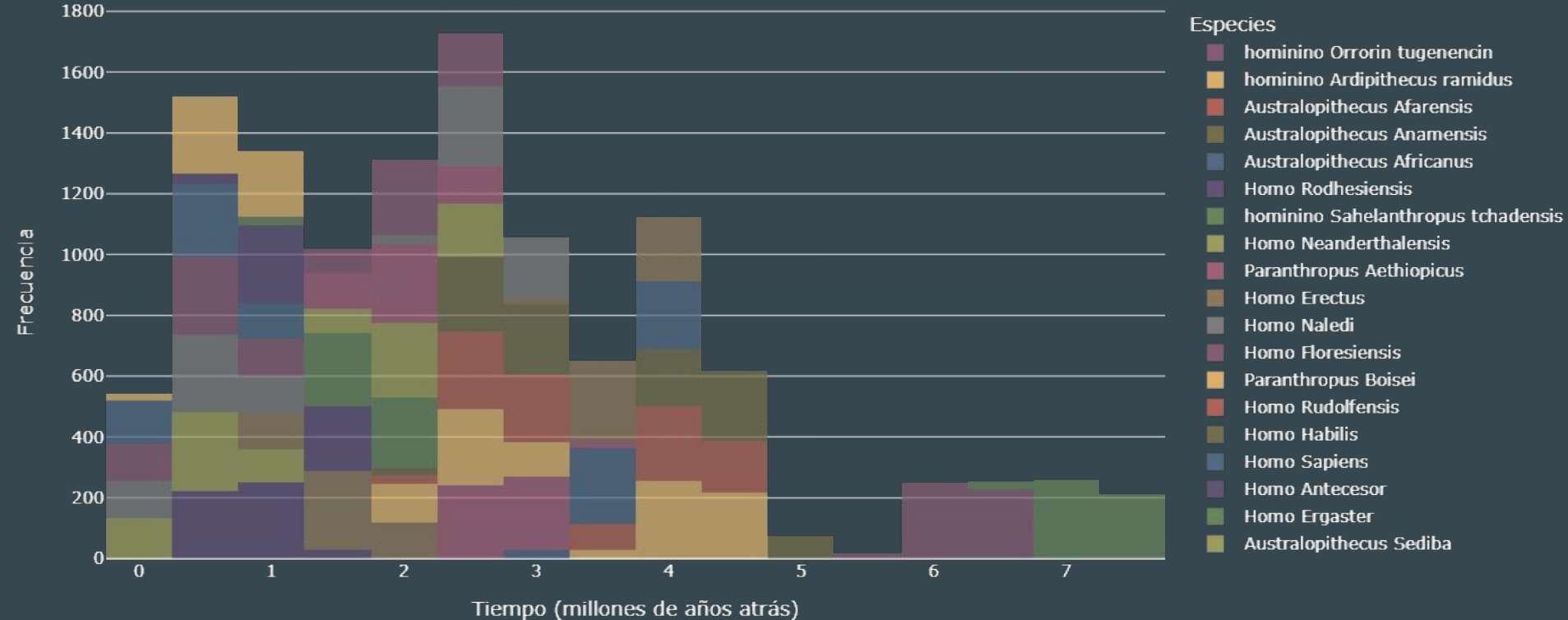
# EDA (Exploratory Data Analysis)

This balance ensures that no single species dominates the analysis, allowing the physical and temporal characteristics of each group to be evaluated equitably.
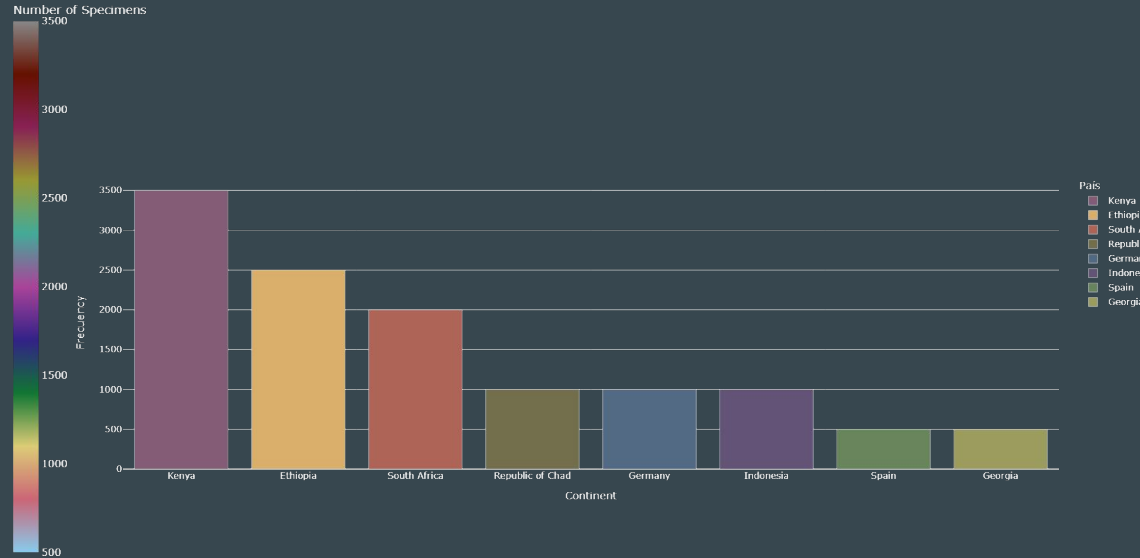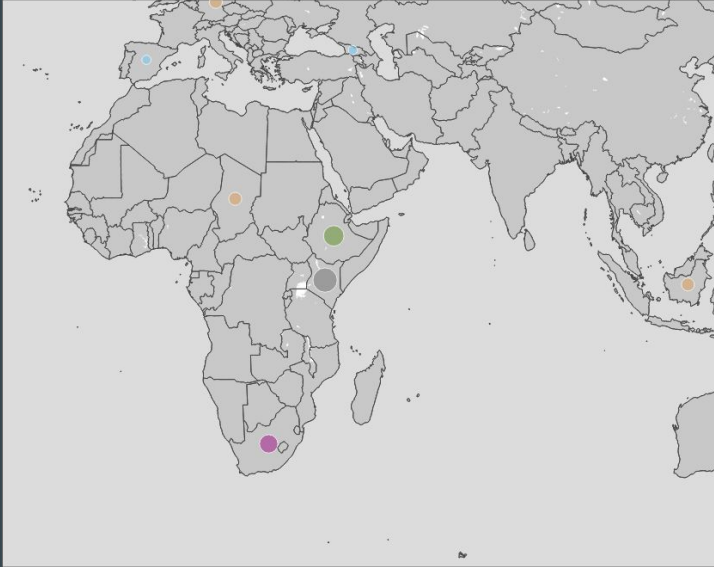


- ■ hominino Orrorin tugenencin
- ■ hominino Ardipithecus ramidus
- ■ Paranthropus Robustus
- ■ Australopithecus Garhi
- ■ Australopithecus Bahrelghazali
- ■ Homo Georgicus
- ■ Australopithecus Sediba
- ■ Homo Ergaster
- ■ Homo Antecesor
- ■ Homo Sapiens
- ■ Homo Habilis
- ■ Homo Rudolfensis
- ■ Paranthropus Boisei
- ■ Homo Floresiensis
- ■ Homo Naledi
- ■ Homo Erectus
- ■ Paranthropus Aethiopicus
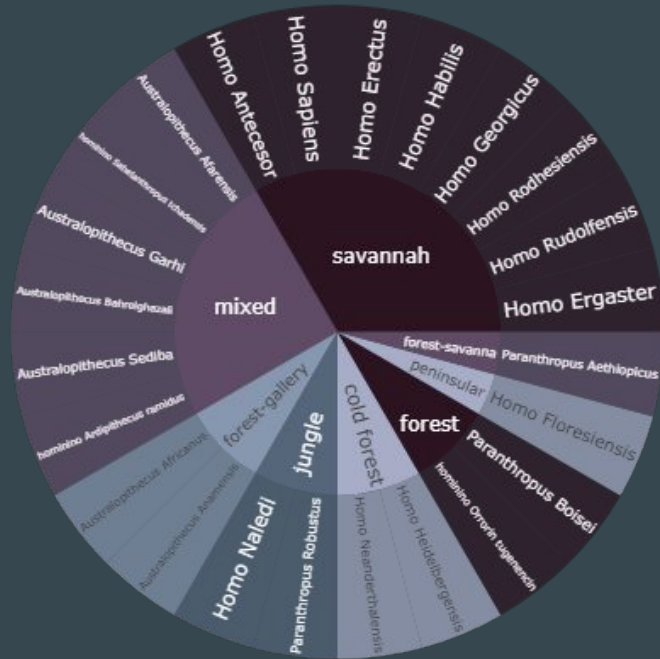- ■ Homo Neanderthalensis

# Species Over Time

Frequency of fossil records in each time interval.

# Number of Fossils by Country
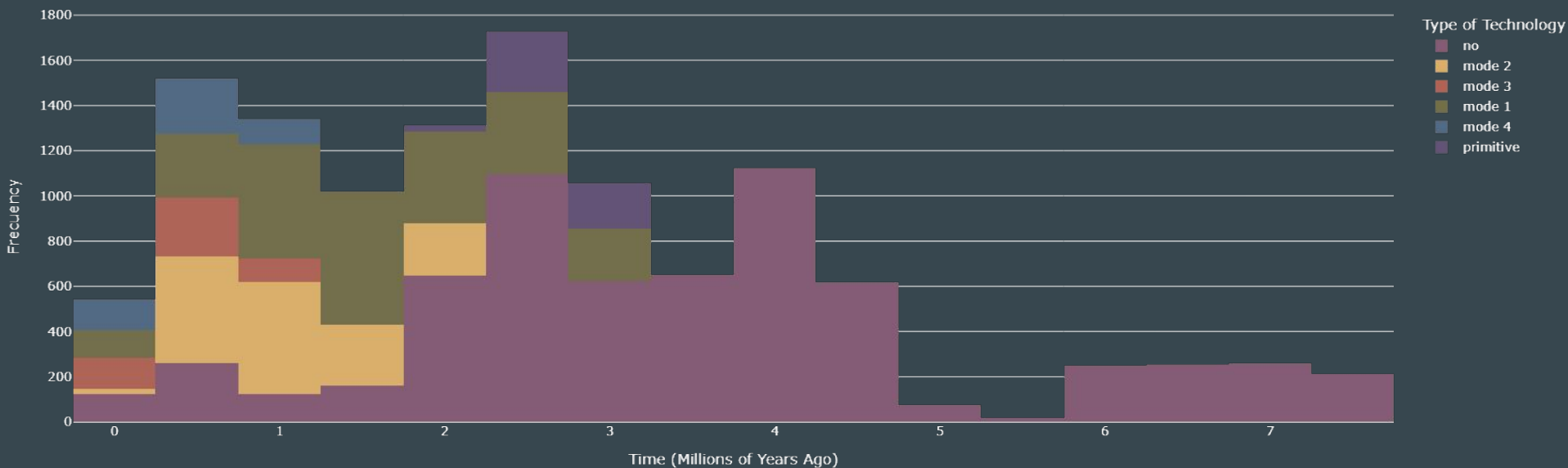
# Habitat by Species



There is a close relationship between hominid species diversity and their ecological environments. The predominance of **savanna habitats** highlights their role as a key evolutionary cradle, closely associated with species such as *Homo erectus*, *Homo habilis*, and *Homo ergaster*.

Meanwhile, more specific habitats reflect how different species adapted to varying environmental pressures. This provides a comprehensive evolutionary perspective on how hominids interacted with and thrived in their surroundings.
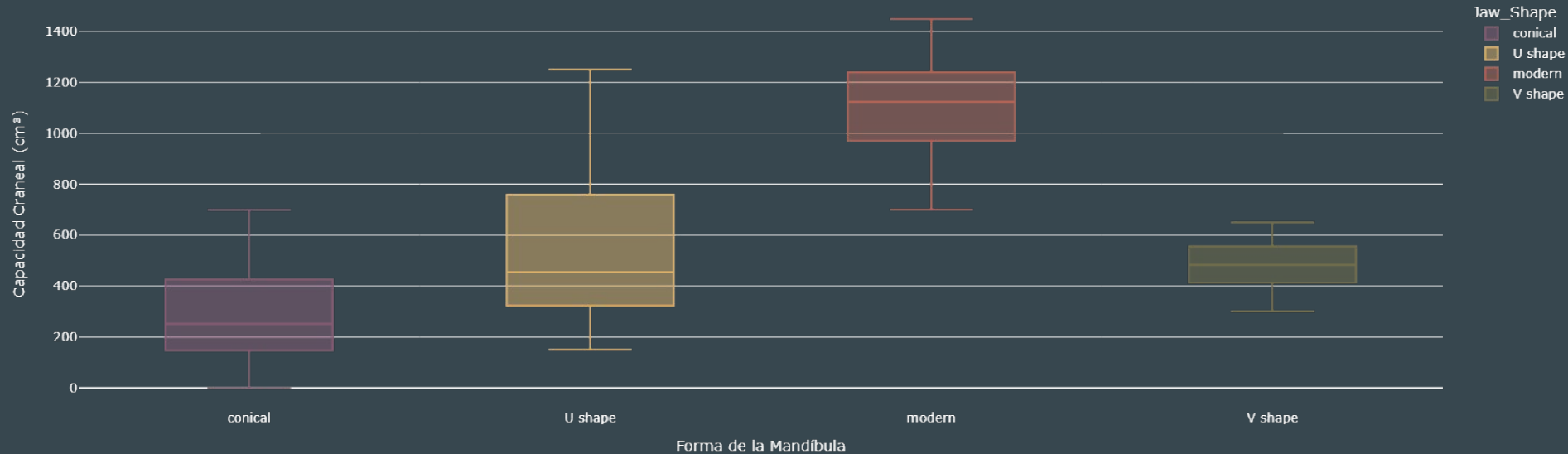
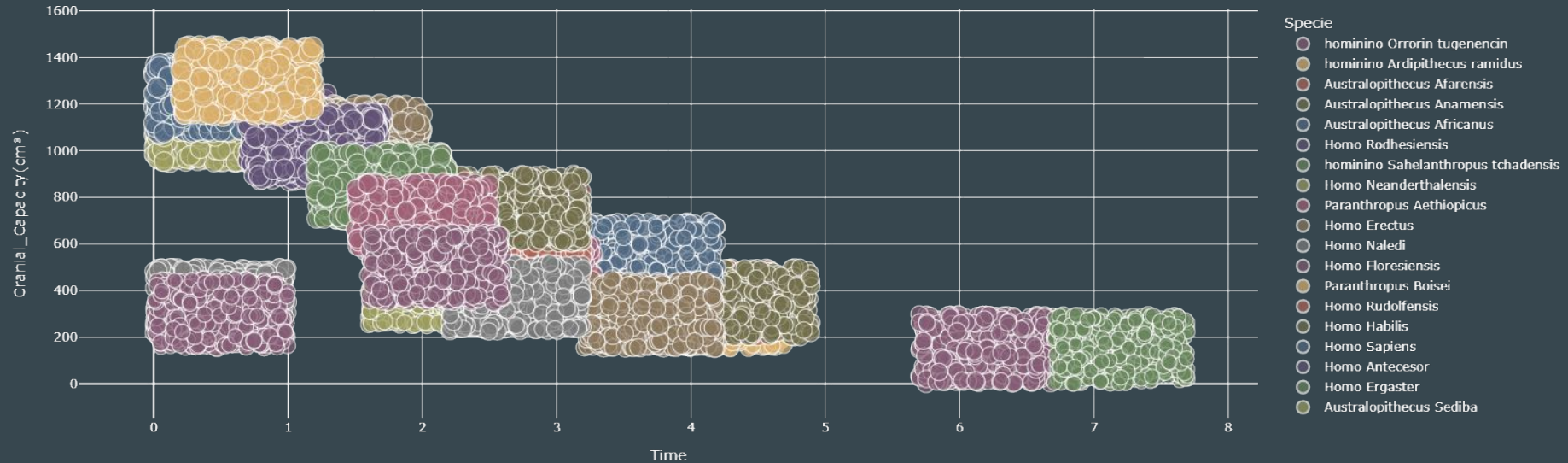Hierarchical Distribution of Tecno_Type and Genus & Species

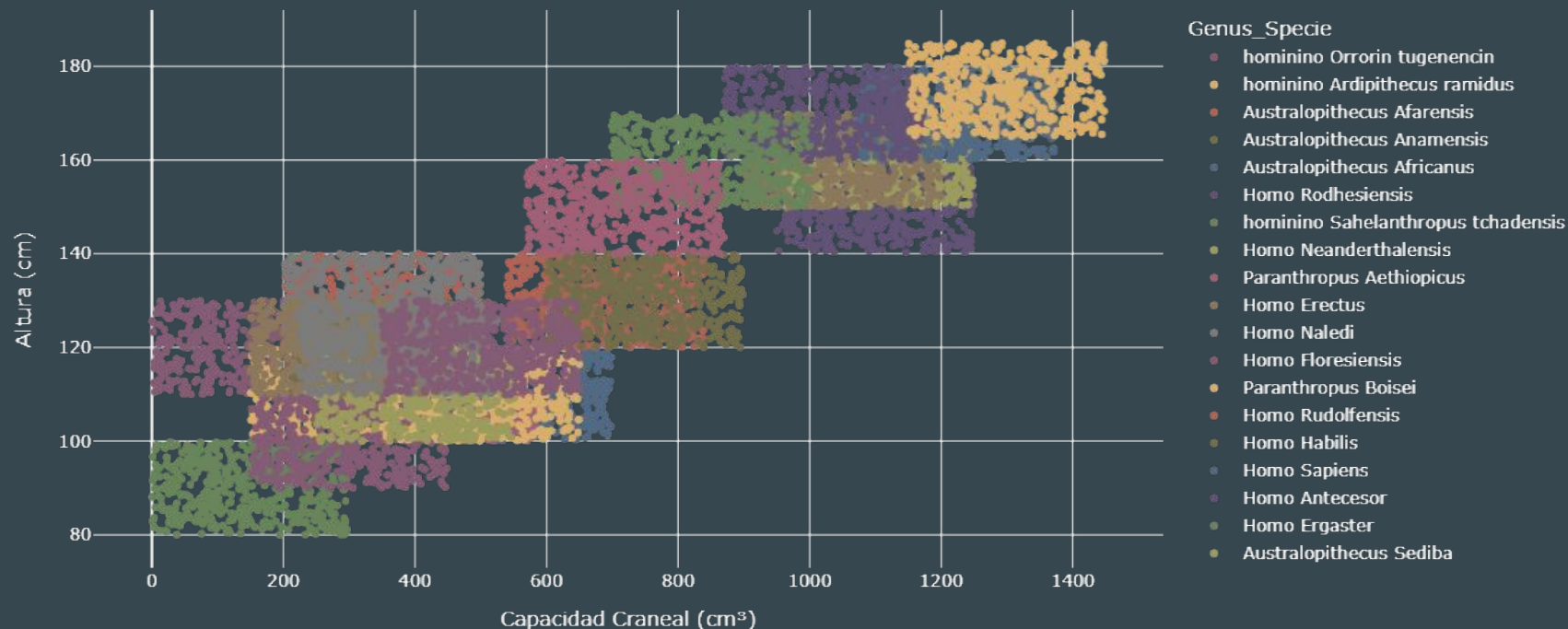Relationship Between Cranial Capacity and Jaw Shape

Relationship Between Geological Time (in Millions of Years) and Cranial Capacity

Evolutionary Trend: A progressive increase in cranial capacity is observed as time decreases (more recent species).

Correlation Between Cranial Capacity and Height

# Correlation of Categorical Variables (Cramér's V)



Correlation Matrix (Cramér's V)

- Certain categorical variables, such as Genus_Species, Habitat, and Diet, show strong correlations with other features:
- Genus_Species: Highly correlated with variables like Tecno, Tecno_type, and Habitat, indicating that these characteristics are closely linked to species classification.
- Diet: Strongly correlated with Canine Size and Tooth Enamel, reflecting specific dietary adaptations.
- Tecno_type: Moderately correlated with Biped, Diet, and Habitat, suggesting its relationship with anatomical and ecological changes.

# Correlation of Numerical Variables



- Very strong positive correlation: This suggests that species with larger cranial capacity also tend to have greater height.

- Moderate negative correlation: More recent species tend to have a larger cranial capacity.

- Moderate negative correlation: More recent species also tend to be taller.

# PCA(Principal Component Analysis)

EThe graph presents a two-dimensional visualization of hominid species based on the first two Principal Components (PC1 and PC2), generated using Principal Component Analysis (PCA).

These components capture the majority of variability present in the selected numerical variables: Time, Cranial Capacity, and Height.

# Data Preprocessing

Scaling of Numerical Variables:

- Min-Max Scaler was used to scale numerical variables since the Shapiro-Wilks test indicated they did not follow a Gaussian distribution.

- This ensured that variable values were within a uniform range (0 to 1), optimizing their usability in the models.

Encoding of Categorical Variables:

- Ordinal Variables: Features with an inherent order, such as Canine Size and Incisor Size, were encoded using Ordinal Encoding, preserving their hierarchical structure.

- Non-Ordinal Variables: Categorical variables without a specific order, such as Habitat, Genus_Species, and Jaw Shape, were encoded using Label Encoding, assigning a unique number to each category without imposing hierarchies.

# Classification Models for Species Prediction

The Random Forest Classifier combines the results of multiple decision trees (through voting) to improve accuracy and reduce overfitting, leveraging the diversity of samples used in its construction.

It was chosen for its effectiveness in handling a mix of numerical and categorical variables, making it well-suited for this dataset.

The Multilayer Perceptron (MLP) is an artificial neural network designed to learn complex nonlinear patterns.

Its architecture includes multiple dense layers with ReLU activations, allowing it to capture deep relationships between the physical characteristics of the fossils.

**01**
Split Train Text
Target: Genus Specie
20% Test
80% TRain

**02**
Training
Random forest Classifier
MLP Classifier
GradientBoostingClassifier

**03**
Hiperparameter
Optimization
RandomSearchCV
Cross Validation

**04**
Evaluation
Precision
Recall
$R^2$

# Principales resultados de Random Forest



Matriz de Confusión - Random Forest



Importancia

The model is classifying most instances correctly, as the majority of values are on the main diagonal.

| precision | recall | f1-score |
|---|---|---|
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |
| 1.0 | 1.0 | 1.0 |

# Resultados de Multilayer Perceptron, MLP



Precisión durante el entrenamiento



Pérdida durante el entrenamiento

```
Reporte de Clasificación:
              precision    recall  f1-score   support

           0       1.00      1.00      1.00       107
           1       1.00      1.00      1.00        92
           2       1.00      1.00      1.00       105
           3       1.00      1.00      1.00       108
           4       1.00      1.00      1.00       103
           5       1.00      1.00      1.00        98
           6       1.00      1.00      1.00        92
           7       1.00      1.00      1.00        90
           8       1.00      1.00      1.00       108
           9       1.00      1.00      1.00        86
          10       1.00      1.00      1.00        97
          11       1.00      1.00      1.00        92
          12       1.00      1.00      1.00        99
          13       1.00      1.00      1.00        85
          14       1.00      1.00      1.00        97
          15       1.00      1.00      1.00        98
          16       1.00      1.00      1.00       107
          17       1.00      1.00      1.00       101
          18       1.00      1.00      1.00       115
          19       1.00      1.00      1.00       105
          20       1.00      1.00      1.00       117
          21       1.00      1.00      1.00        93
          22       1.00      1.00      1.00       113
          23       1.00      1.00      1.00        92

    accuracy                           1.00      2400
   macro avg       1.00      1.00      1.00      2400
weighted avg       1.00      1.00      1.00      2400
```
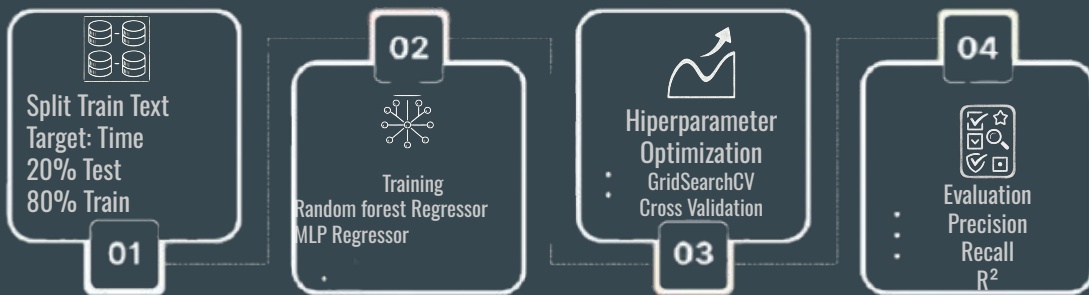
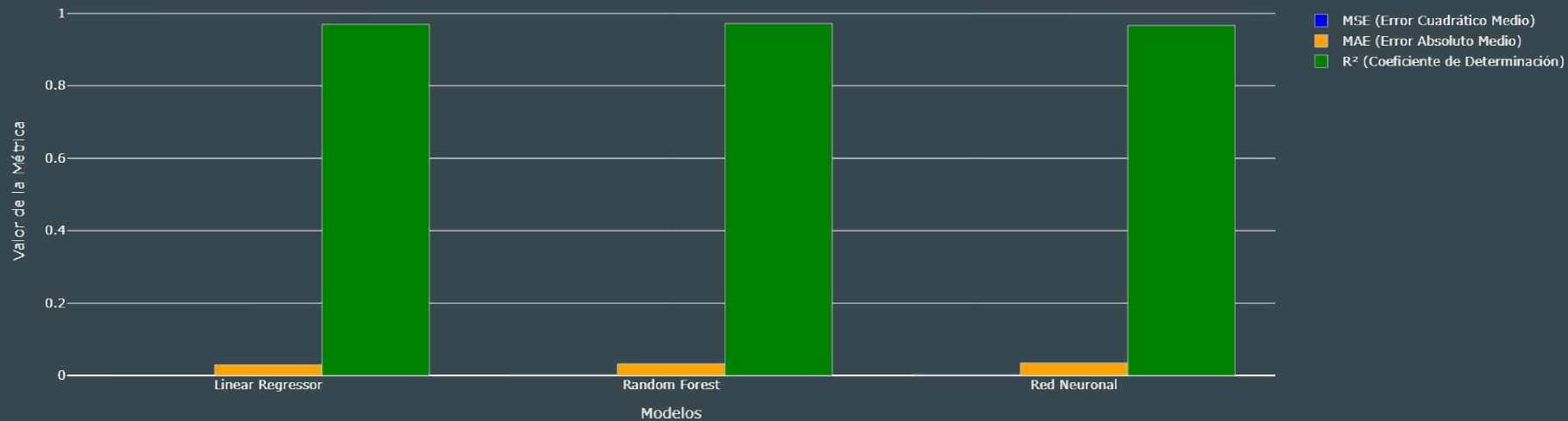❊ The model is achieving excellent performance in both training and validation.

❊ The low loss values in both sets suggest that the model is learning correctly and has a good fit.

# Modelos de regresión para datación preliminar



01 — Split Train Text
Target: Time
20% Test
80% Train

02 — Training
Random forest Regressor
MLP Regressor

03 — Hiperparameter
Optimization
GridSearchCV
Cross Validation

04 — Evaluation
Precision
Recall
$R^2$

|  | MAE | MSE | $R^2$ |
|---|---|---|---|
| LinearRegression | 0.03 | 0.00 | 0.97 |
| Random Forest Regressor | 0.0328 | 0.0328 | 0.9720 |
| Perceptrón Multicapa (MLP) | 0.03505 | 0.00173 | 0.9667 |

# Conclusions

The development and implementation of the models demonstrated the feasibility of using data science for taxonomic classification.

The use of physical characteristics to predict genera and species suggests that these variables contain consistent patterns that can be leveraged for classification purposes. This enables faster and more accurate taxonomic identification.

The models have proven effective in estimating geological time with minimal errors, which can be useful in contexts where absolute dating is unavailable.

The choice of model depends on the objective:

Taxonomic classification: Random Forest.

Time estimation: Linear regression for simple cases, neural networks or Random Forest for more complex relationships.

# Limitations and Future Work

Limitations

Subjectivity in Measurement: Some anatomical characteristics can be difficult to quantify objectively, introducing a certain degree of subjectivity in the data.

Taxonomic Criteria: The definition of species and genus in paleontology is a complex and evolving subject. Classification criteria may vary among researchers, making it challenging to compare results.

Future Work

Expanding the Dataset: Incorporating new fossil samples, particularly those with uncommon traits or from understudied geographic regions.

Establishing a Reproducible Framework: Creating a methodology that can be replicated for future data-driven studies. This approach could also be transferred to forensic anthropology for predicting characteristics in modern individuals.