

## Feature Engineering

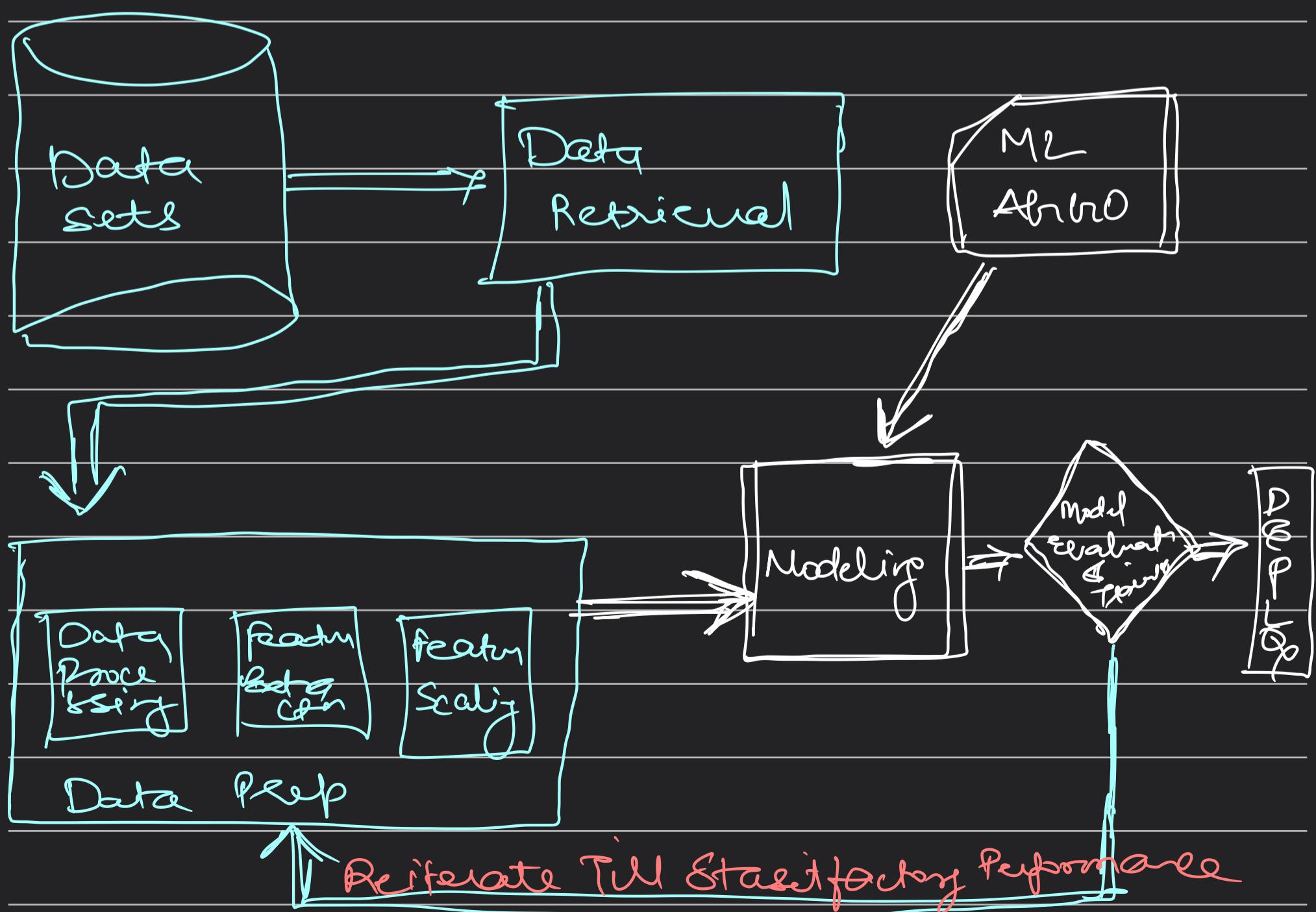
↳ Is the process of using domain knowledge to extract features from raw data.

Then these features can be utilized to improve performance of ML Algorithms.

It's more of an ART → Feature Engineering.

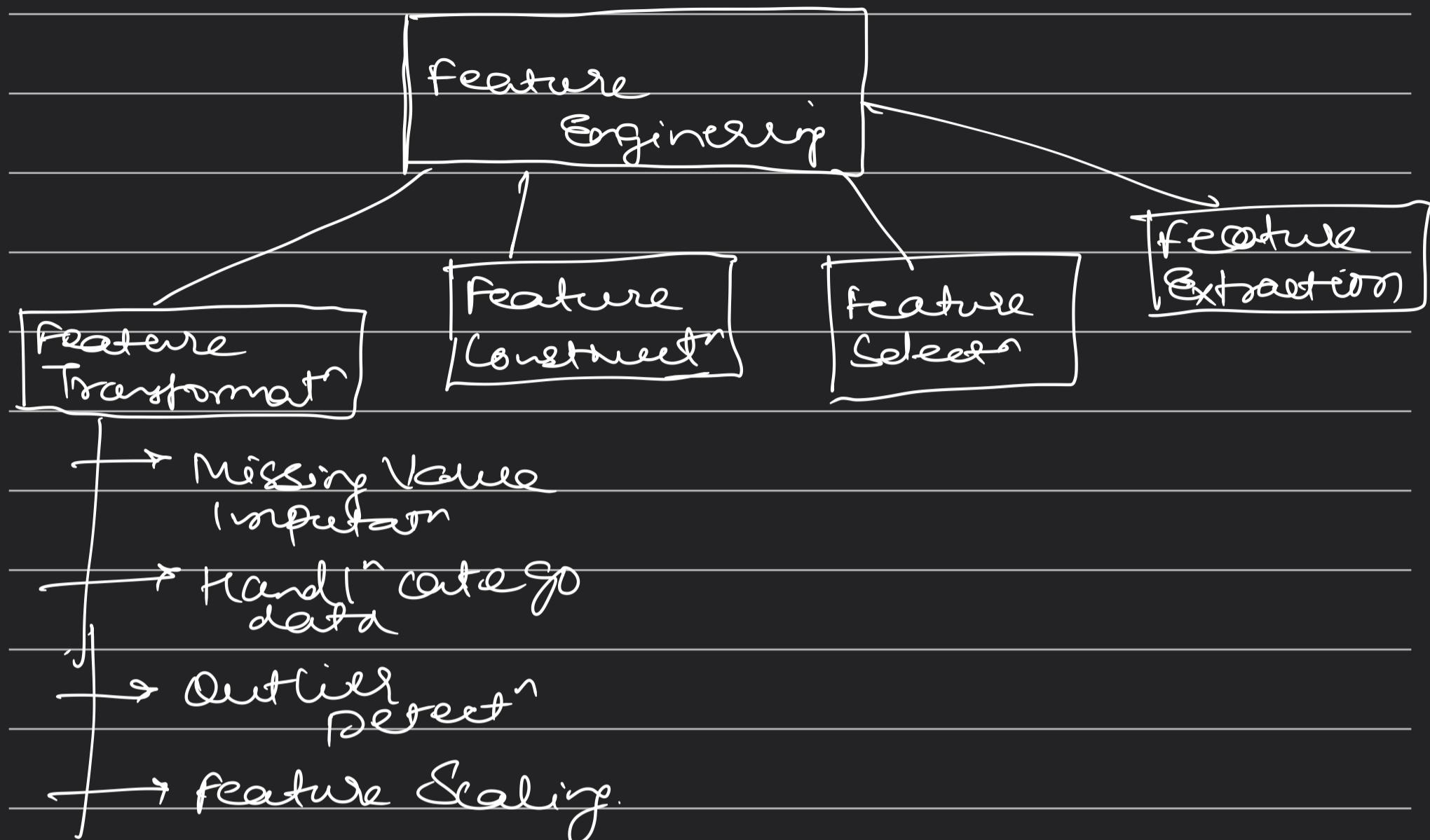
"Even a bad algorithm will perform extremely well if its features are fine grained Picket"

"On the other hand a good algo will perform bad if it doesn't get good features"



ML flow for any Problem.

# Flow chart of feature Engineering



## Feature Transformation

→ Handling Missing Values



This first step in feature Engineering

- Handling by adding values
- Handling by removing values.

→ Handling Categorical Value

- Your libraries only understand numerical data. So your categories must get converted into some numerical value.

→ Ex - One hot encoding.

→ Outlier Detection

Outliers are noisy data, they cause the

Algo to behave abnormally when trained & give drastic result.

### → Feature Scaling

Few features will have dominance in value over others, this can impact our algo. For this we scale features.

Ex- Min Max Scaling, Normalization etc.

Ex Age & Salary =  $\frac{32}{2}, \frac{100000}{1}$  \$  
This difference in  $\Leftrightarrow$  **Big Difference**  
Scaling of features effect the Algo.

### 2 Feature Construction

We create a complete new column or variable or feature based upon our intuition we get from domain knowledge of data.

→ Family Member column can be created using Sibling & Par Arch column in Titanic Database.

→ We can also create categorical data called Alone with 0 family Members & Family Man etc.

### 3 Feature Selection

We have lot of input columns & we select based upon our domain knowledge

which variables or features we are going to use in our Machine learning algorithm.

Now we start to operate on less features because we eliminated the unnecessary ones.

### Feature Extraction

Programmatically we create new features based upon the features provided.

It's bit different from feature construction

For Ex  $\Rightarrow$  Room, WashRoom, Price These are necessary data for Real Estate Predictor but we can use Room + Wash Room = Carpet Area.

This way new feature got created & we got rid of the other two.

In a way we reduce features & use the new created ones.

Used specially in high Dimensional Data.

