

Project Title: **Forecasting House Prices Accurately Using Smart Regression Techniques in Data Science**

GITHUB REPOSITORY LINK: <https://github.com/Vaithisan12/House.git>

1. Problem Statement:

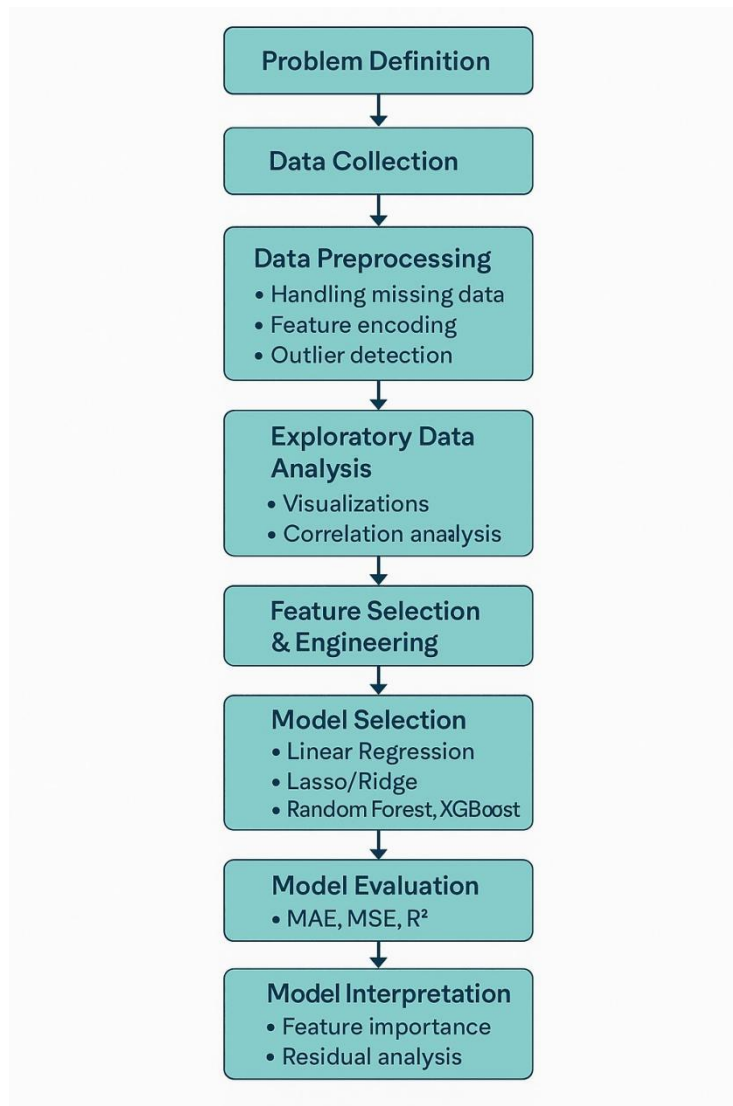
Accurate prediction of house prices is a critical challenge in the real estate industry, where property values are influenced by numerous factors such as location, size, number of bedrooms, age of the house, and market trends. Traditional pricing models often fall short in capturing complex, non-linear relationships between these features and housing prices. The objective of this project is to develop a robust and intelligent regression-based model using advanced data science techniques to forecast house prices with high accuracy. By leveraging smart regression algorithms—such as regularized linear models (Lasso, Ridge), ensemble methods (Random Forest, Gradient Boosting), and advanced approaches like XGBoost or neural networks—the goal is to build a model that can generalize well across diverse housing datasets.

2. Project Objectives:

1. **To collect and preprocess real-world housing datasets** containing relevant features such as location, square footage, number of rooms, age, and neighborhood attributes.
2. **To perform exploratory data analysis (EDA)** to identify patterns, outliers, correlations, and feature importance that influence house prices.
3. **To implement and compare multiple smart regression techniques** such as:
 - Linear Regression
 - Ridge and Lasso Regression
 - Decision Trees
 - Random Forest Regression
 - Gradient Boosting (including XGBoost, LightGBM)
4. **To optimize model performance** using hyperparameter tuning, cross-validation, and regularization techniques to avoid overfitting and improve generalization.

5. **To evaluate model accuracy** using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and R-squared score (R^2).

3. Flowchart of the project workflow



4. Data Description

The dataset used in this project contains various features (independent variables) that influence house prices (the dependent variable). Below is a general description of the key attributes typically found in housing price datasets:

Feature Name	Description
Id	Unique identifier for each house entry
SalePrice	Final selling price of the house (Target variable)
LotArea	Lot size in square feet
OverallQual	Overall material and finish quality (rated 1–10)
OverallCond	Overall condition rating (rated 1–10)
YearBuilt	Year the house was built
YearRemodAdd	Year of last remodeling
TotalBsmtSF	Total basement area in square feet
GrLivArea	Above ground (gross living area) square feet
GarageArea	Size of the garage in square feet
GarageCars	Number of cars that can fit in the garage
FullBath	Number of full bathrooms
HalfBath	Number of half bathrooms
BedroomAbvGr	Number of bedrooms above ground level
KitchenQual	Kitchen quality (e.g., Ex, Gd, TA, Fa)
Neighborhood	Physical locations within the city (categorical feature)
HouseStyle	Style of dwelling (e.g., 1Story, 2Story)
MSZoning	General zoning classification (e.g., RL, RM)
Exterior1st	Exterior covering on house (e.g., VinylSd, MetalSd)
Heating	Type of heating system
CentralAir	Whether the house has central air conditioning (Y/N)

5. Data Preprocessing

1. Importing Required Libraries & Loading Data

- Use libraries like `pandas`, `numpy`, `matplotlib`, and `seaborn` to load and explore the dataset.

2. Handling Missing Values

- **Numerical Features:** Impute with mean/median.
- **Categorical Features:** Impute with mode or a new category (e.g., 'None').

- Drop columns with too many missing values if necessary.

3. Encoding Categorical Variables

- **Label Encoding:** For ordinal data (e.g., quality ratings).
- **One-Hot Encoding:** For nominal categories (e.g., neighborhoods, house styles).

4. Outlier Detection and Removal

- Use box plots or Z-scores to detect outliers in features like `GrLivArea`, `SalePrice`, etc.
- Consider removing or capping extreme outliers to avoid model distortion.

5. Feature Engineering

- Create new features from existing ones, such as:
 - `AgeOfHouse = YrSold - YearBuilt`
 - `TotalBathrooms = FullBath + 0.5 * HalfBath`
 - `TotalSF = TotalBsmtSF + GrLivArea + GarageArea`

6. Exploratory Data Analysis (EDA)

EDA helps understand the structure, patterns, and relationships within the housing dataset before building any model. Here are the key steps:

1. Data Overview

- Display basic info using `df.info()` and `df.describe()` to understand:
 - Data types (numeric, categorical)
 - Summary statistics
 - Missing values

2. Target Variable Analysis

- **Visualize `SalePrice` (target)**
 - Histogram or KDE Plot: `sns.histplot(df['SalePrice'], kde=True)`
 - Check distribution (often right-skewed)
 - Apply log transform if highly skewed

3. Correlation Analysis

- **Correlation matrix:** `df.corr()`
 - **Heatmap:** Use `sns.heatmap()` to find features strongly correlated with `SalePrice`
 - Common high-correlation features: `GrLivArea`, `OverallQual`, `GarageArea`, `TotalBsmtSF`
-

4. Univariate Analysis

- **Numerical Features:**
 - Histograms, boxplots to identify distribution and outliers
 - Example: `sns.boxplot(x=df['GrLivArea'])`
 - **Categorical Features:**
 - Count plots: `sns.countplot(x='HouseStyle', data=df)`
 - Bar plots: Compare mean `SalePrice` across categories
-

5. Bivariate Analysis

- **Numeric vs Target:**
 - Scatterplots: `sns.scatterplot(x='GrLivArea', y='SalePrice', data=df)`
 - Linearity and outliers can be checked
- **Categorical vs Target:**
 - Boxplots: `sns.boxplot(x='Neighborhood', y='SalePrice', data=df)`
 - Helps identify high-value area

7. Feature Engineering

Feature engineering improves model accuracy by creating, transforming, or selecting the most informative variables. It's a critical step in housing price prediction due to the complexity of real estate data.

1. Creating New Features

1. Total Square Footage:

python:

```
df['TotalSF'] = df['TotalBsmtSF'] + df['1stFlrSF'] + df['2ndFlrSF']
```

2. Categorical Variables

- **Label Encoding** (for ordinal features):

python:

```
from sklearn.preprocessing import LabelEncoder
```

```
df['KitchenQual'] = LabelEncoder().fit_transform(df['KitchenQual'])
```

3. Handling Skewed Features

- Apply log transformation to reduce skew:

python:

```
df['SalePrice'] = np.log1p(df['SalePrice'])
```

```
df['GrLivArea'] = np.log1p(df['GrLivArea'])
```

4. Feature Selection

- Drop redundant or highly correlated features:

python:

```
df = df.drop(['GarageCars'], axis=1) # If strongly correlated with GarageArea
```

5. Polynomial Features (Optional)

- Add interaction or squared terms for linear models:

python:

```
from sklearn.preprocessing import PolynomialFeatures
```

```
poly = PolynomialFeatures(degree=2, include_bias=False)
```

```
X_poly = poly.fit_transform(df[['GrLivArea', 'TotalSF']])
```

8. Model Building Steps

1. Import Required Libraries

Python:

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
```

```
from sklearn.linear_model import LinearRegression, Ridge, Lasso
```

```
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
```

```
from xgboost import XGBRegressor
```

```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
```

2. Model Interpretation

- **Feature Importance:**

python:

```
import matplotlib.pyplot as plt
```

```
feat_importances = pd.Series(rf.feature_importances_, index=X.columns)
```

```
feat_importances.nlargest(10).plot(kind='barh')
```

```
plt.title("Top 10 Feature Importances")
```

```
plt.show()
```

9. isualization of results & Model Insights

1. Actual vs Predicted Plot

Visualizes how well the model is predicting compared to real values.

Python:

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
y_pred = model.predict(X_test)
```

```
plt.figure(figsize=(8,6))
```

```
sns.scatterplot(x=y_test, y=y_pred)
```

```
plt.xlabel('Actual Sale Price')
```

```
plt.ylabel('Predicted Sale Price')
```

```
plt.title('Actual vs Predicted House Prices')
```

```
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', lw=2)
```

```
plt.show()
```

10. Tools and Technologies Used:

1. Programming Languages:

- **Python:** The most widely used language for data science, with libraries like pandas, NumPy, scikit-learn, and statsmodels for regression and analysis.
- **R:** Another strong option for statistical modeling, with packages like `caret`, `xgboost`, and `randomForest`.

2. Data Science Libraries:

- **pandas:** For data manipulation and cleaning.
- **NumPy:** For numerical computations.
- **scikit-learn:** For machine learning algorithms, including regression models and model evaluation tools.
- **statsmodels:** For more advanced statistical models like ordinary least squares regression.
- **XGBoost, LightGBM, CatBoost:** Libraries for powerful boosting algorithms.
- **Matplotlib, Seaborn:** For data visualization.

3. Deep Learning (Optional):

- **TensorFlow or PyTorch:** For more complex models like neural networks (if the dataset is large and requires deep learning techniques).

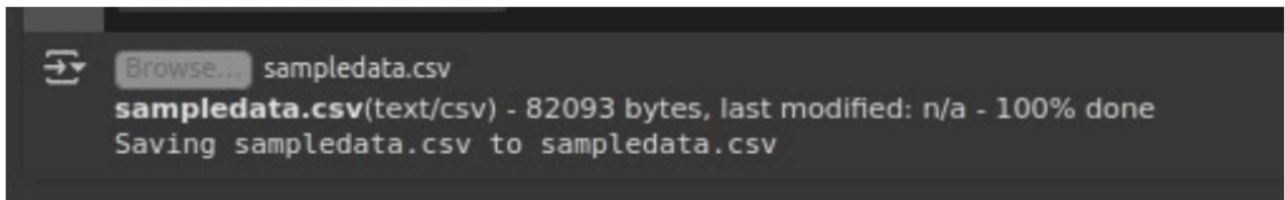
4. Data Handling:

- **SQL:** To query databases and retrieve housing data.
- **NoSQL (MongoDB, Firebase):** For unstructured or semi-structured data.

5. Cloud Platforms (for Deployment and Scaling):

- **AWS** (Amazon Web Services), **Google Cloud**, or **Microsoft Azure**: For hosting models, using managed machine learning services, and storing large datasets.

```
from google.colab import files
uploaded=files.upload()
```

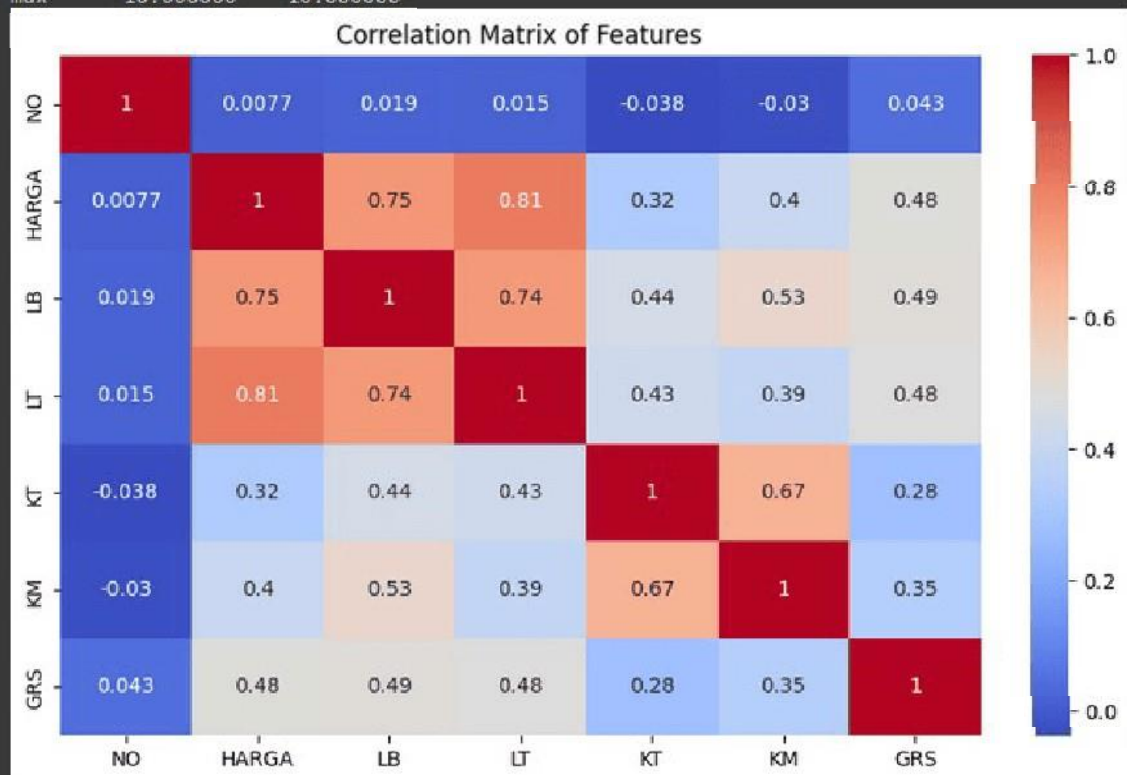


```
import matplotlib.pyplot as plt
import seaborn as sns
# Basic statistics
print(df.describe())
# Correlation matrix
plt.figure(figsize=(10, 6))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title("Correlation Matrix of Features")
plt.show()
```



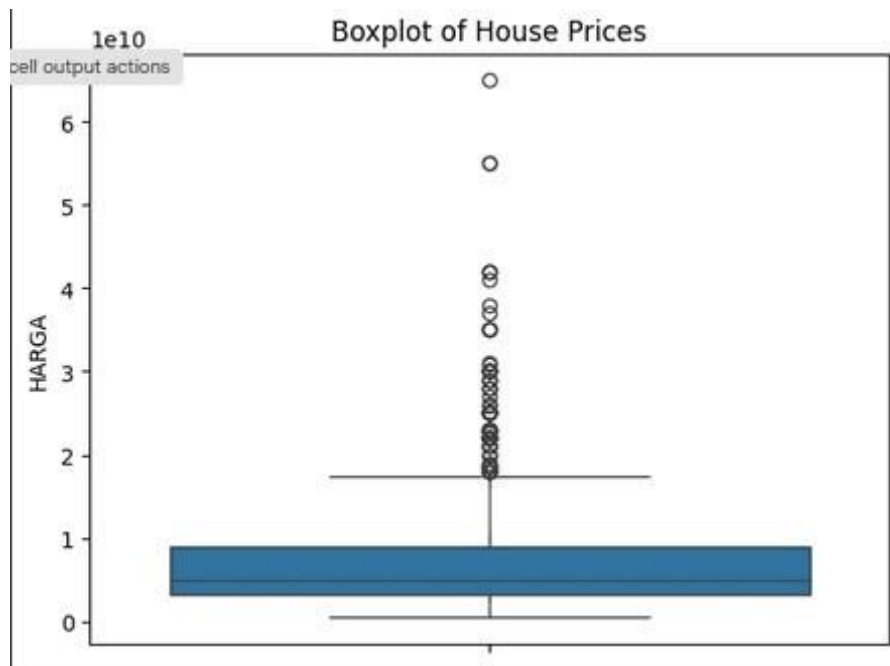
	NO	HARGA	LB	LT	KT
count	1010.000000	1.010000e+03	1010.000000	1010.000000	1010.000000
mean	505.500000	7.628987e+09	276.539604	237.432673	4.668317
std	291.706188	7.340946e+09	177.864557	179.957604	1.572776
min	1.000000	4.300000e+08	40.000000	25.000000	2.000000
25%	253.250000	3.262500e+09	150.000000	130.000000	4.000000
50%	505.500000	5.000000e+09	216.500000	165.000000	4.000000
75%	757.750000	9.000000e+09	350.000000	290.000000	5.000000
max	1010.000000	6.500000e+10	1126.000000	1400.000000	10.000000

	KM	GRS
count	1010.000000	1010.000000
mean	3.607921	1.920792
std	1.420066	1.510998
min	1.000000	0.000000
25%	3.000000	1.000000
50%	3.000000	2.000000
75%	4.000000	2.000000
max	10.000000	10.000000



```
# Boxplot to detect price outliers
sns.boxplot(df['HARGA'])
plt.title("Boxplot of House Prices")
plt.show()
```

```
# Remove extreme outliers
df_clean = df[df['HARGA'] < df['HARGA'].quantile(0.95)]
```



```
from sklearn.linear_model import LinearRegression
import numpy as np
```

```
X = df[['LB']]
y = df['HARGA']
```

```
model = LinearRegression()
model.fit(X, y)
```

```
print("R^2 Score:", model.score(X, y))
```

```
R^2 Score: 0.5581327856561413
```

```
features = ['LB', 'LT', 'KT', 'KM', 'GRS']
X = df[features]
y = df['HARGA']
```

```
model = LinearRegression()
model.fit(X, y)
```

```
print("Model coefficients:", model.coef_)
print("R^2 score:", model.score(X, y))
```

```
Model coefficients: [ 1.23187516e+07  2.36590867e+07 -6.19514797e+08  4.55486747e+08
 3.09965160e+08]
R^2 score: 0.7162361438094645
```

```
from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline

poly_model = make_pipeline(PolynomialFeatures(2), LinearRegression())
poly_model.fit(X, y)

print("R^2 score (poly):", poly_model.score(X, y))
```

```
↗ R^2 score (poly): 0.7384637627019939
```

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
model = LinearRegression()
model.fit(X_train, y_train)
pred = model.predict(X_test)

print("Test RMSE:", np.sqrt(mean_squared_error(y_test, pred)))
```

```
↗ Test RMSE: 2986616943.9349093
```

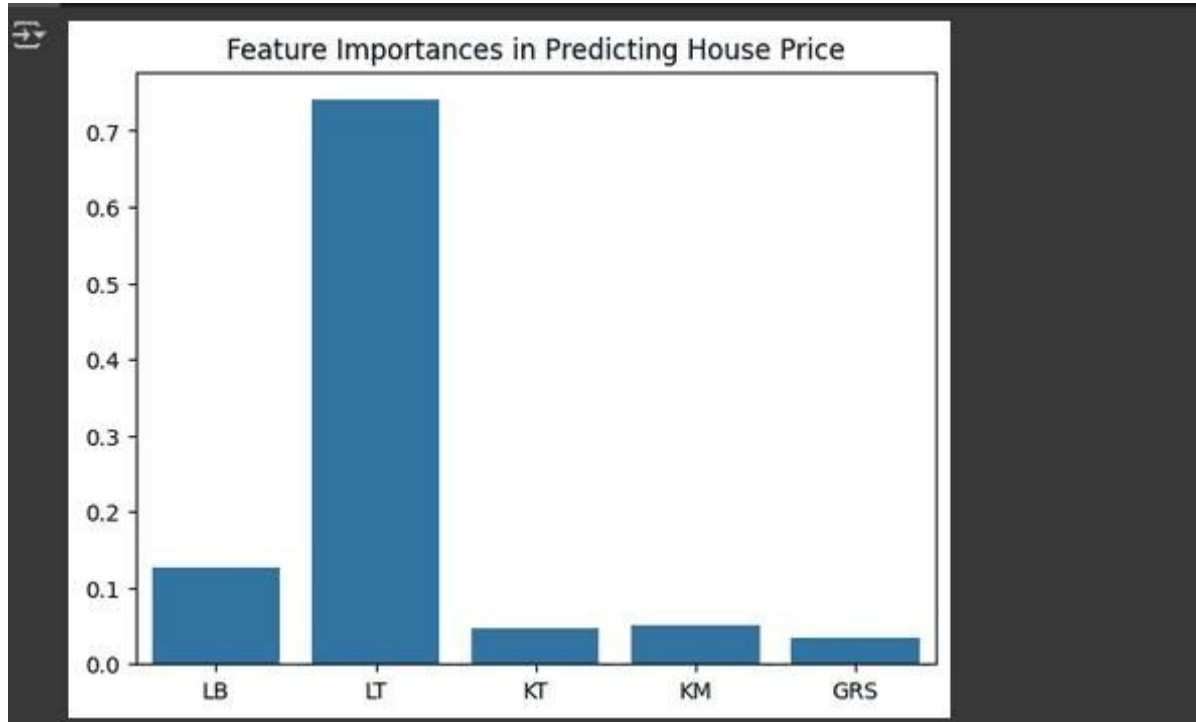
```
from sklearn.ensemble import RandomForestRegressor

rf = RandomForestRegressor(n_estimators=100)
rf.fit(X_train, y_train)

print("Random Forest R^2:", rf.score(X_test, y_test))
```

```
↗ Random Forest R^2: 0.8498985772208579
```

```
importances = rf.feature_importances_  
sns.barplot(x=features, y=importances)  
plt.title("Feature Importances in Predicting House Price")  
plt.show()
```



```
sample = pd.DataFrame({  
    'LB': [200],  
    'LT': [150],  
    'KT': [4],  
    'KM': [3],  
    'GRS': [1]  
})
```

```
predicted_price = rf.predict(sample)  
print("Predicted House Price:", predicted_price[0])
```

```
Predicted House Price: 4644500000.0
```

```
import joblib

joblib.dump(rf, 'house_price_model.pkl')
print("Model saved as 'house_price_model.pkl'")
```

A screenshot of a terminal window with a dark background. It shows the output of a Python command: "Model saved as 'house_price_model.pkl'". There is a small icon on the left side of the terminal window.

Team Members and Contributions

1.SURENDER-PROBLEM STATEMENT AND PROJECT OBJECTIVES

2.SRIDHAR-FLOWCHART OF THE PROJECT WORKFLOW,DATA DESCRIPTION,DATA PREPROCESSING,EXPLORATORYDATA ANALYSIS

3.SUNIL-FEATURE ENGINEERING AND MODEL BULIDING

4.VATHISH- VISULATION OF RESULTS&MODEL INSIGHTS AND TOOLS&TECHNOLOGIES USED