# Analysis of Camera Dataset using Machine Learning Algorithms

TEAM DATA PIRATES

Pawan Prasad P, Phani Kumar Vedurumudi, G U Deepak, Jayanth O

PES UNIVERSITY EC CAMPUS

Hosur Rd, Konappana Agrahara, Electronic City, Bengaluru, Karnataka-560100

## ABSTRACT

*Camera is used in every day to day life to capture the wonderful memories of our life. A camera is an optical instrument that captures a visual image. As we know that the evolution of camera from the ancient time to the recent days camera has a dramatical changes in the features like lens, aperture, focal length, price etc. We have taken the dataset to check the radical changes in the price over a period from 1997-2002. There are many varieties of cameras from olden days camera to Digital camera. Our main focus is on the price prediction of camera based on features of camera models using Machine Learning Algorithms like Linear Regression, Lasso-Ridge Regression, KNN and SVR Models. We basically used Regression Models because our dataset has continuous data. The features with more correlation are used for analysis. The overall Accuracy rate or prediction rate can reach upto 70% but generally it is between 50 – 70%, because the camera dataset doesn't have much correlation among the features.*
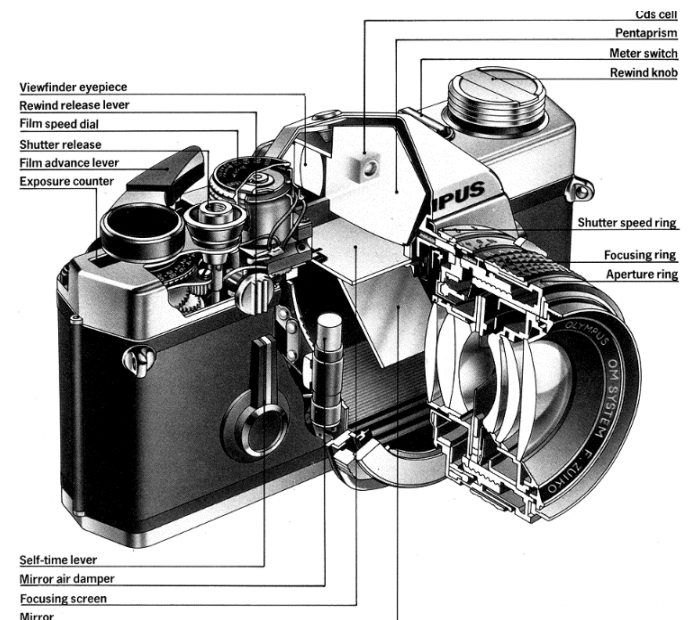
*Keywords—Camera, features, continuous data Period, price prediction, Machine Learning Algorithms, Linear Regression, Lasso-Ridge Regression, KNN, SVR, Accuracy rate.*

## I. INTRODUCTION

A camera is an optical instrument that captures a visual image. At a basic level, cameras are sealed boxes (camera body) with a small hole (aperture) that allows light through to capture a image on light-sensitive surface (usually photographic film or a digital sensor).

Our Dataset has different kinds of camera models with their features like Model, Release date, Max resolution, Low resolution, Effective pixels, Zoom wide (W), Zoom tele (T), Normal focus range, Macro focus range, Storage included, Weight (inc. batteries), Dimensions and Price. The Model feature is having categorical data and all other features are having continuous data. Due to **continuous data** we have used Regression Models to predict the price of the models. Our dataset has 13 features and 1038 observations.



The main focus of our project is to predict the price of camera based on features of camera models using Machine Learning Algorithms like Linear Regression, Lasso-Ridge Regression, KNN and SVR Models. We took problem related to price prediction because we have nowadays share market predictions which is a great problem in this tech world.

## II. PROBLEM STATEMENT

Our Project objective is to predict the price of the model based on the given features of the camera model. We have used the various Machine Learning Algorithms Models like KNN (k-nearest neighbours), Linear Regression and SVR.

## III. RELATED WORK

We have done the literature survey of nearly 12 recent journal Papers in which the authors have used various theory to use the camera in various AI fields like Computer Vision and AI in Security. But we have used different kinds of models to do the price prediction which is our main motive and we have chosen the price prediction which is a regression problem and due to more no. of types of models we cannot find the model type with all other features of camera which consists continuous data (value). Due to continuous data of all the features we have selected various ML Regression Models related to our dataset which is more precise and gives the correct prediction of price. The author's of many journal Papers have tried only one ML model for their project but we have used the best model out of various Regression models like SVR (Support Vector Regressor, KNN (k-nearest neighbours), Linear Regression Model and Lasso-Ridge Regression Model. The author of all this journal papers have used Neural Network, Human Tracking Algorithm, Camera Calibration Algorithm and many more. We have took all the useful techniques from the previous works. Our main motive is to find the prediction of prices of new data and our dataset have only the continuous value so we have took the various Regression Models to solve our problem. The scope of our problem is basically to give a clear cut prediction of the price of camera models and use the best model to deal with this problem.
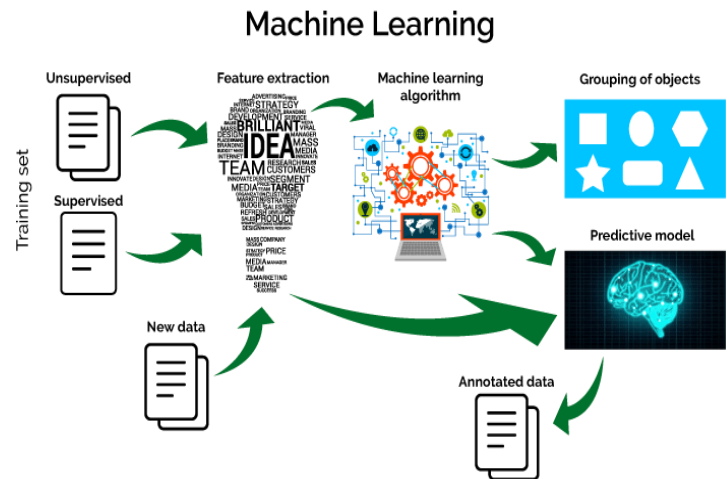
## IV. MODEL BUILDING AND TESTING

The Model Building Process is very crucial part and it is done in a systematic manner. We have used various Supervised Machine Learning Models for the price prediction and the tool which we have used for all the process of model building and testing is the **scikit learn** which is a very good tool for ML.



The Model is build using the following steps:

- ❖ Loading the dataset from the CSV File
- ❖ Exploratory Data Analysis (EDA)
- ❖ Preprocessing the Data
- ❖ Data Visualization
- ❖ Building Regression Model
- ❖ Cross Validation (Testing and Training (80-20 split) and Performance Metrics
- ❖ Model Evaluation
- ❖ Comparison Analysis



❖ **Loading the dataset from the CSV File**

We have used the basic and advance libraries of the python programming language to build our model. Some of them are numpy, pandas, matplotlib, seaborn and sklearn etc.

Our Data has been collected from Kaggle's dataset which is online and free dataset. 1000 Cameras Dataset have been gathered and cleaned up by Petra Isenberg, Pierre Dragicevic and Yvonne Jansen.

```
import math
import pandas as pd
import matplotlib.pyplot  as plt
import numpy as np
import sklearn
import sklearn.preprocessing as pre
import sklearn.linear_model as lm
import seaborn as sns
import sklearn.datasets
import sklearn.neighbors as nb
import sklearn.pipeline as pipeline
from sklearn.preprocessing import StandardScaler
import sklearn.svm as svm
import sklearn.neural_network as nn
import sklearn.metrics as metrics
import sklearn.tree as tree
from sklearn.model_selection import cross_val_score
```

```
df=pd.read_csv('E:\PES COLLEGE\SEM 5\Data Analytics Project\camera_dataset.csv')
```

❖ **Exploratory Data Analysis (EDA)**

This Dataset contains all the details of different Camera Models features and their values which helps to analysis various real-world problems. It has both Categorical and continuous data.

1. Dataset name: Camera Dataset
2. Source: Kaggle
3. No. of Observations: 1038
4. No. of Columns: 13
5. Percentage of Missing Values: 5%

The no. of Missing Values in the Dataset.



This Dataset contains a combination of numerical and categorical features:

Categorical:
- Model

Numerical:
- Release date
- Max resolution
- Low resolution
- Effective pixels

- Zoom wide (W)
- Zoom tele (T)
- Normal focus range
- Macro focus range
- Storage included
- Weight (inc. batteries)
- Dimensions
- Price

Removing the Missing or Incomplete values in the dataset.



❖ **Preprocessing the Data**

The data of the dataset shows less correlation among the attributes. So we have selected only those attributes which has more correlation for model building and we have done random sampling for changing the data randomly for training and testing. We have used Normalization and Standardization which done in the Pipeline of sklearn.
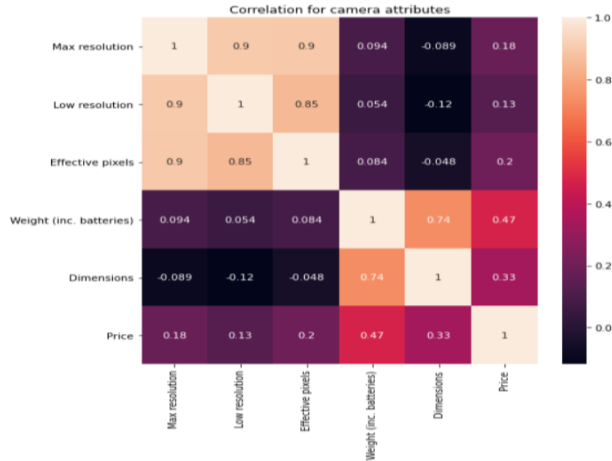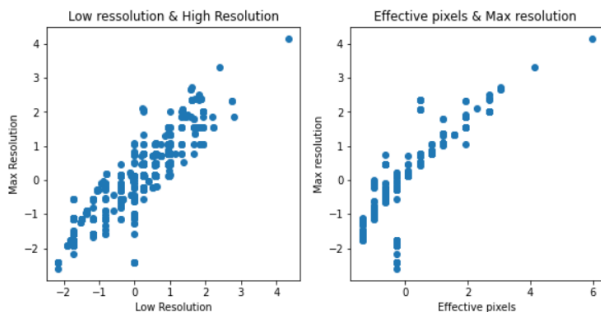
### ❖ Data Visualization

**Heat map:**

The Heat map shows the correlation between the features. These features are used for the model building and testing.



Correlation for camera attributes

**Scatter Plot:**

This plot shows the distribution of the data in the dataset. The variation or distance from each point to another point.



### ❖ Building Regression Model and Cross Validation and Evaluation

We have built many models and I would like to address only two models among them that is the Linear Regression Model and KNN Model. We have did the 5-fold Cross Validation which is like 80-20 split of data. The Evaluation of Model is done using various Performance Metrics like RMSE, MAE and Rsquare. The Model is built using sklearn by creating a Pipeline which is a model and which would preprocess the data using Standard scalar. We have used Performance Metrics also to evaluate the Model. Please check out our Github link below.

**Linear Regression Model and Cross validation:**



**KNN Model and Cross Validation:**



### ❖ Comparison Analysis or Result

From the RMSE, MSE and Rsquare scores of the above all the Regression models we can make a inference that KNN model is better than all other models.

RMSE of Linear Regression - 0.8756
R^2 of Lasso and Ridge Regression - 0.1269
RMSE of KNN - 0.7293
RMSE of SVR - 0.9456

But we can clearly see that RMSE value of KNN Model is less compare to other models. In this Model evaluation we have consider the normalized data and we came to conclusion that KNN is the best model for prediction.

### V. CONCLUSION

The Dataset which we have used consist of less correlation among the features so our accuracy is around 50 to 70% which is given by KNN Model. We have not split the data explicitly to test the model instead we have used the 5-fold cross validation test which automatically splits the data in 80-20 split and based various Performance Metrics scores we

have came to an conclusion that KNN Model is best for our dataset. We have also performed the tuning of parameters of the models by eliminating the less correlated features (attributes) from the dataset and then trained the model. The other Regression models are not so good for our dataset and we can conclude that KNN Model is best for our problem that is Price Prediction of the camera model.

## VI. REFERENCES

[1] https://scikit-learn.org/stable/
[2] https://scikit-learn.org/stable/modules/cross_validation.html
[3] https://www.jstor.org/stable/25750706
[4] https://ieeexplore.ieee.org/abstract/document/709612/
[5] https://www.sciencedirect.com/science/article/pii/S0379073820304631

### Entire Project Details:

❖ https://github.com/Pawantech-App/DATA-ANALYTICS-PROJECT-

❖ https://drive.google.com/drive/folders/1kb9kgaIb6HjDvNBAdP683aatBqJ9bX1M?usp=sharing

## CONTRIBUTION:

Team Members:

• Pawan Prasad P - PES2UG19CS280
• Phani KumarVedurumudi  - PES2UG19CS281
• G U Deepak - PES2UG19CS124
• Jayanth O - PES2UG19CS163