

IE 7275: Data Mining in Engineering
Housing Price Prediction

Milestone: Project Report

Group11

Gaddam Sreeramulu Rohit Kumar

Mobile Student : (857)-265-0719

gaddamsreeramulu.r@northeastern.edu

Signature of Student : Rohit Kumar

Submission Date: 04-28-2023

Title of the Project: House Pricing Prediction

Problem Setting

The concept of house pricing is the process of determining the value of a house based on various factors such as location, size, condition, amenities, and market demand. House pricing is a critical component of the real estate industry and is used to set the asking price for a house when it is put up for sale.

Challenges:

1. Higher interest rates have stopped buyers in their tracks.
2. House inventory has gotten low.
3. Getting the right price for the house is tough.

Application:

There are several applications of house pricing models in the real world, including:

Real estate valuation: House pricing models can be used to estimate the value of a property for sales or rental purposes.

Homeownership: Homeowners can use house pricing models to understand the market value of their property, to make decisions about when to sell or make improvements, and to plan for their financial future.

Urban planning and development: Urban planners and developers can use house pricing models to make informed decisions about land use, zoning, and investment opportunities.

Problem Definition

The problem of predicting the sale price of a house based on the various features like size of the house, number of bedrooms, location, etc., can help real estate agents and homeowners in pricing their houses correctly, and can also assist potential buyers in determining a fair price for a house they are interested in. The target variable (price) can be influenced by these factors and the goal is to build a model that accurately predicts the price of a house based on these factors.

Data Source

In this project for house pricing prediction, we will be using data set from 'Kaggle', which is a Google LLC subsidiary, where users can search data sets and build models in a web-based data science environment.

Here's a snippet of our dataset after performing the data cleaning and preprocessing, where we removed null and duplicate values and dropped unnecessary columns.

Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
79545.45857431680	5.682861321615590	7.009188142792240	4.09	23086.800502686500	1059033.5578701200	208 Michael Ferry Apt. 674 Laurabury, NE 37010-5101
79248.64245482570	6.0028998082752400	6.730821019094920	3.09	40173.07217364480	1505890.91484695	188 Johnson Views Suite 079 Lake Kathleen, CA 48958
61287.067178656800	5.865889840310000	8.512727430375100	5.13	36882.15939970460	1058987.9878760800	9127 Elizabeth Stravenue Danieltown, WI 06482-3489
63345.24004622800	7.1882360945186400	5.586728664827650	3.26	34310.24283090710	1260616.8066294500	USS Barnett FPO AP 44820
59982.197225708000	5.040554523106280	7.839387785120490	4.23	26354.109472103100	630943.4893385400	USNS Raymond FPO AE 09386
80175.7541594853	4.9884077575337100	6.104512439428880	4.04	26748.428424689700	1068138.0743935300	06039 Jennifer Islands Apt. 443 Tracyport, KS 16077
64698.46342788770	6.025335906887150	8.147759585023430	3.41	60828.24908540720	1502055.8173744100	4759 Daniel Shoals Suite 442 Nguyenburgh, CO 20247
78394.33927753090	6.9897797477182800	6.620477995185030	2.42	36516.358972493800	1573936.5644777200	972 Joyce Viaduct Lake William, TN 17778-6483
59927.66081334960	5.36212556960358	6.3931209805509000	2.3	29387.39600281590	798869.5328331630	USS Gilbert FPO AA 20957
81885.92718409570	4.423671789897880	8.167688003472350	6.1	40149.96574921340	1545154.8126419600	Unit 9446 Box 0958 DPO AE 97025
80527.47208292290	8.09351268063935	5.042746799645980	4.1	47224.35984022190	1707045.722158060	6368 John Motorway Suite 700 Janetbury, NM 26854
50593.69549704280	4.496512793097040	7.467627404008020	4.49	34343.991885578800	663732.3968963270	911 Castillo Park Apt. 717 Davisborough, PW 78603
39033.809236982400	7.671755372854430	7.250029317273500	3.1	39220.36146737250	1042814.0978200900	209 Natasha Stream Suite 961 Huffmanland, NE 52457
73163.6634410467	6.919534825456560	5.9931879009455700	2.27	32326.123139488100	1291331.5184858200	829 Welch Track Apt. 992 North John, AR 26532-5136
69391.3801843616	5.344776176735730	8.406417714534250	4.37	35521.294033173200	1402818.2101658500	PSC 5330, Box 4420 APO AP 08302
79004.86871500000	5.44045040805470	8.517510744107000	4.04	39000.504050000000	1000074.8500510000	8078 Charles View...

Data Description

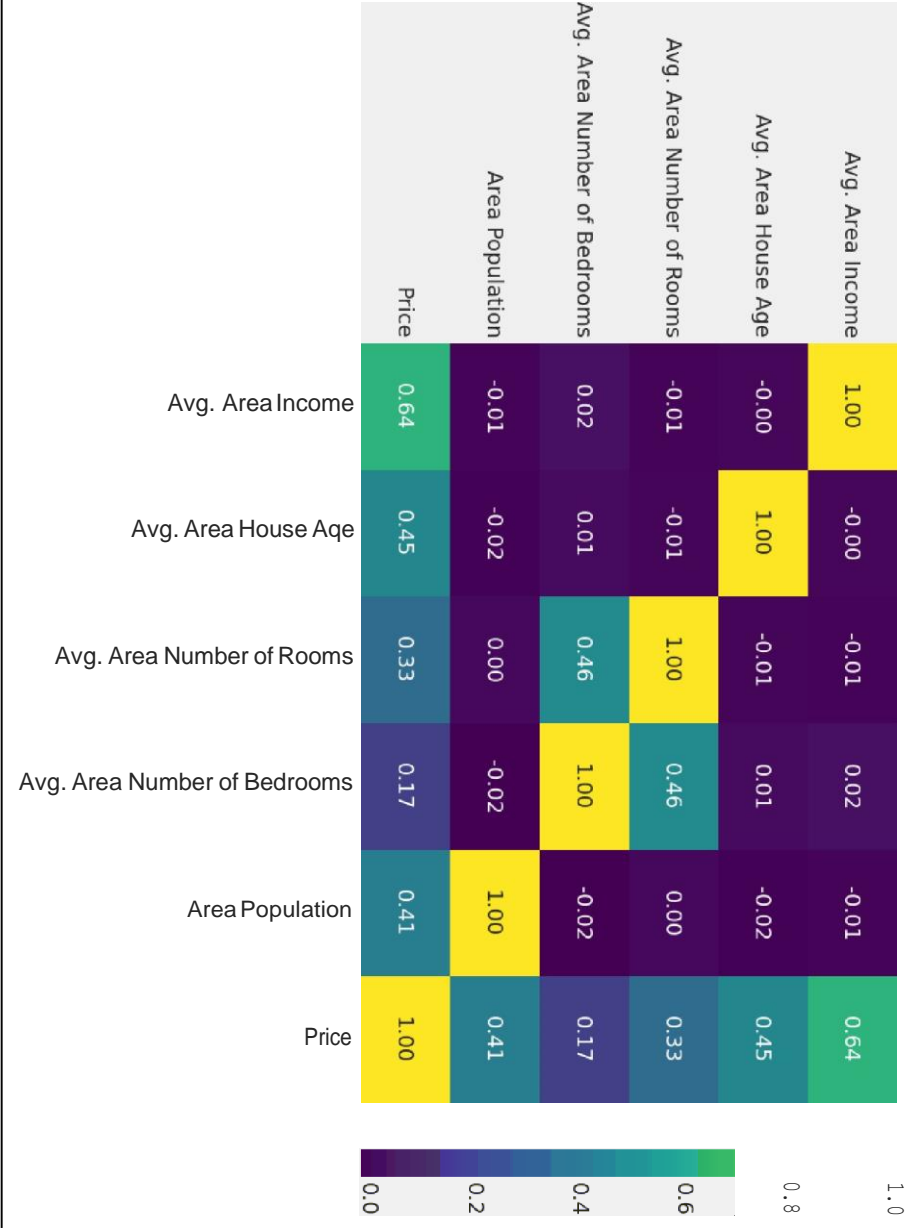
The dataset contains features like sales price, neighborhood, number of bedrooms, number kitchens, locality, sale type, sale condition and so on.

The target variable (price) can be influenced by these factors and the goal is to build a model that accurately predicts the price of a house based on these factors.

We will be using the USA_Housing dataset, which contains information about houses. The objective of this task is to predict the price of a house, which is a continuous variable, making this a regression problem. The dataset includes several columns such as the average income of residents in the city where the house is located, the average age and number of rooms and bedrooms in houses in the same city, the population of the city, the selling price of the house, and the house's address.

Below you can see the heatmaps of our dataset.





Model Exploration and Model Selection

Here, we will try to explore various models and to find the model which has proven to be the most efficient one.

We will explore around with several models before we make the selection.

There are two types of models:

- Predictive Models
- Descriptive Models

Predictive models functions have the order as Classification, Regression, Time Series Analysis, Prediction.

Descriptive models process in the order as Clustering, Summarization, Associative Rules, Sequence.

There are several models such as logistic regression, SVM, KNN, and so on...

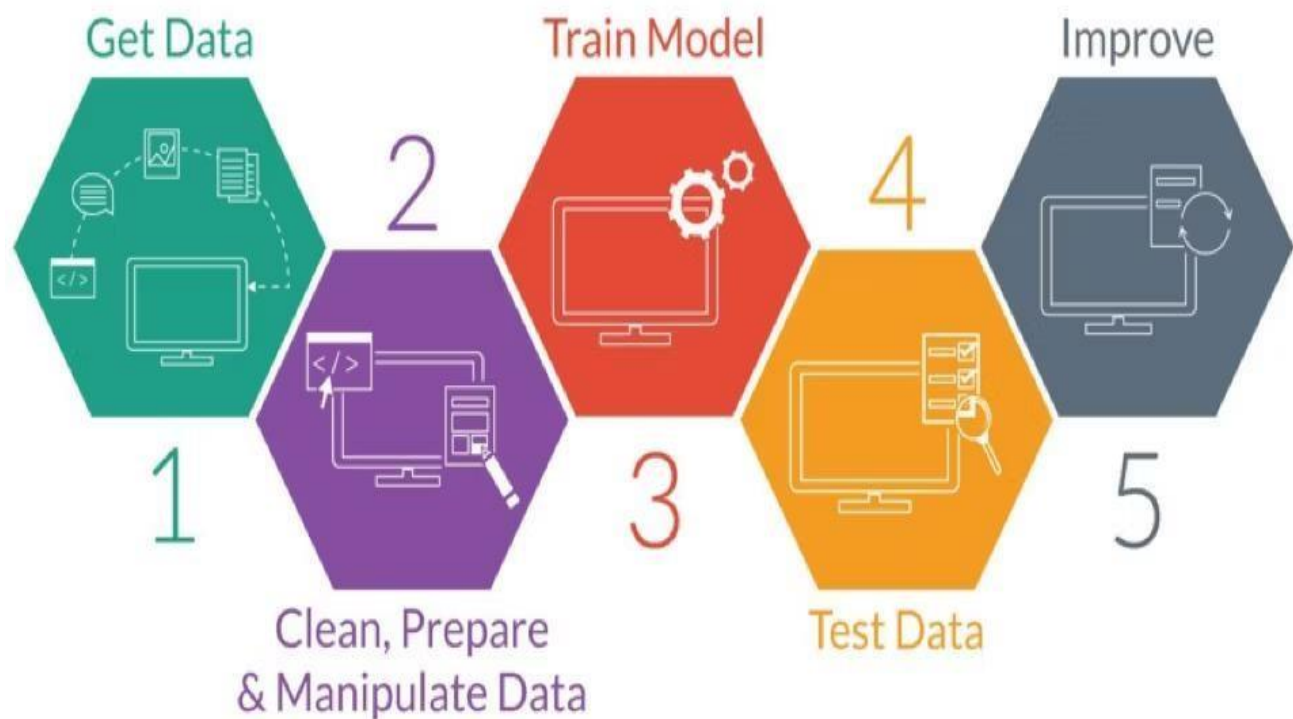
We will make the decision of the model selection by using some of the points as below:

- a model that satisfies the demands and limitations of project stakeholders.
- a model that, given the time and resources at hand, is suitably skilled.
- a model as that is skillful as compared to naive models.
- a model that performs well compared to other models that have been examined.
- a model that is proficiently relative to the state-of-the-art.

In a perfect world, the data would be divided into training, validation, and test sets. Candidate models would then be fitted on the training set, evaluated, and chosen on the validation set, and the performance of the chosen model would be reported on the test set.

Process Flowchart

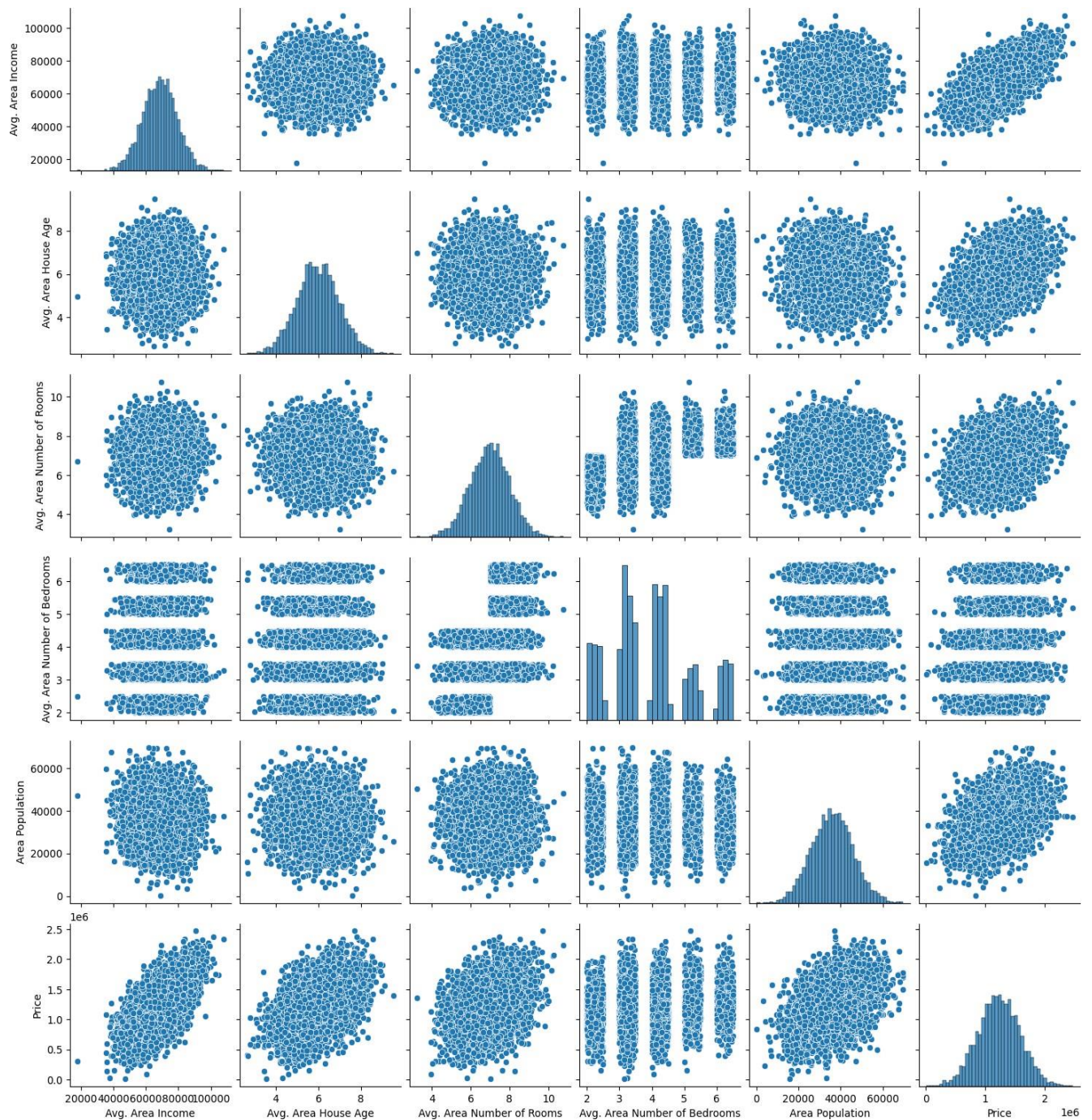
Our project process involves several steps that are essential for extracting useful insights from a large dataset. Firstly, data needs to be collected and curated from various sources. Once the data is collected, it needs to be cleaned and manipulated to remove any errors, missing values, or inconsistencies. After that, machine learning models are trained using the pre-processed data. The model is then tested on a separate set of data to evaluate its performance. Based on the results of the testing, the model can be refined and improved to enhance its accuracy and efficiency. The iterative process of improving the model can be repeated until the desired level of performance is achieved. Overall, this process is an essential step towards unlocking the value of big data and gaining insights that can drive business decisions.



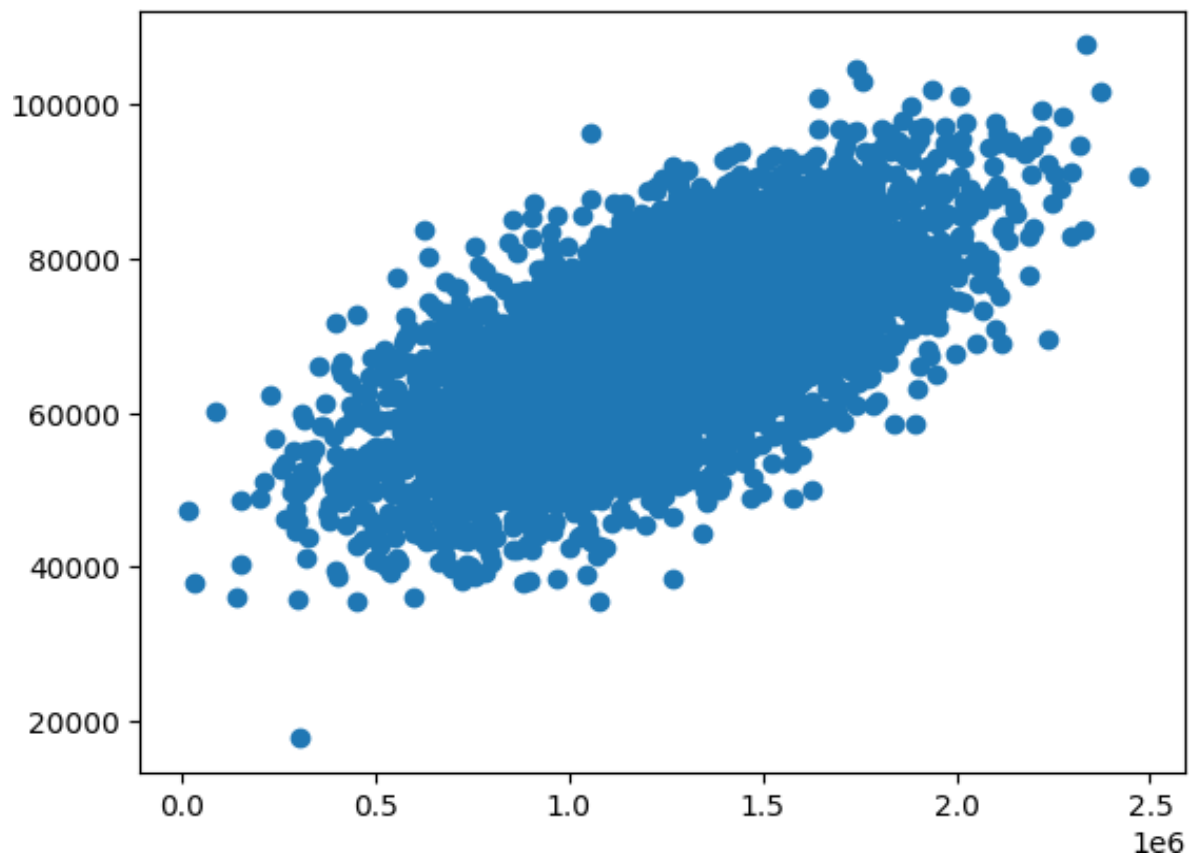
Project Process Flowchart

Data Exploration and Visualization

From our dataset, we have plotted a Pair plot each column in our dataset as seen below such as Avg. Area income, Area population, Avg. House Age, and so on.

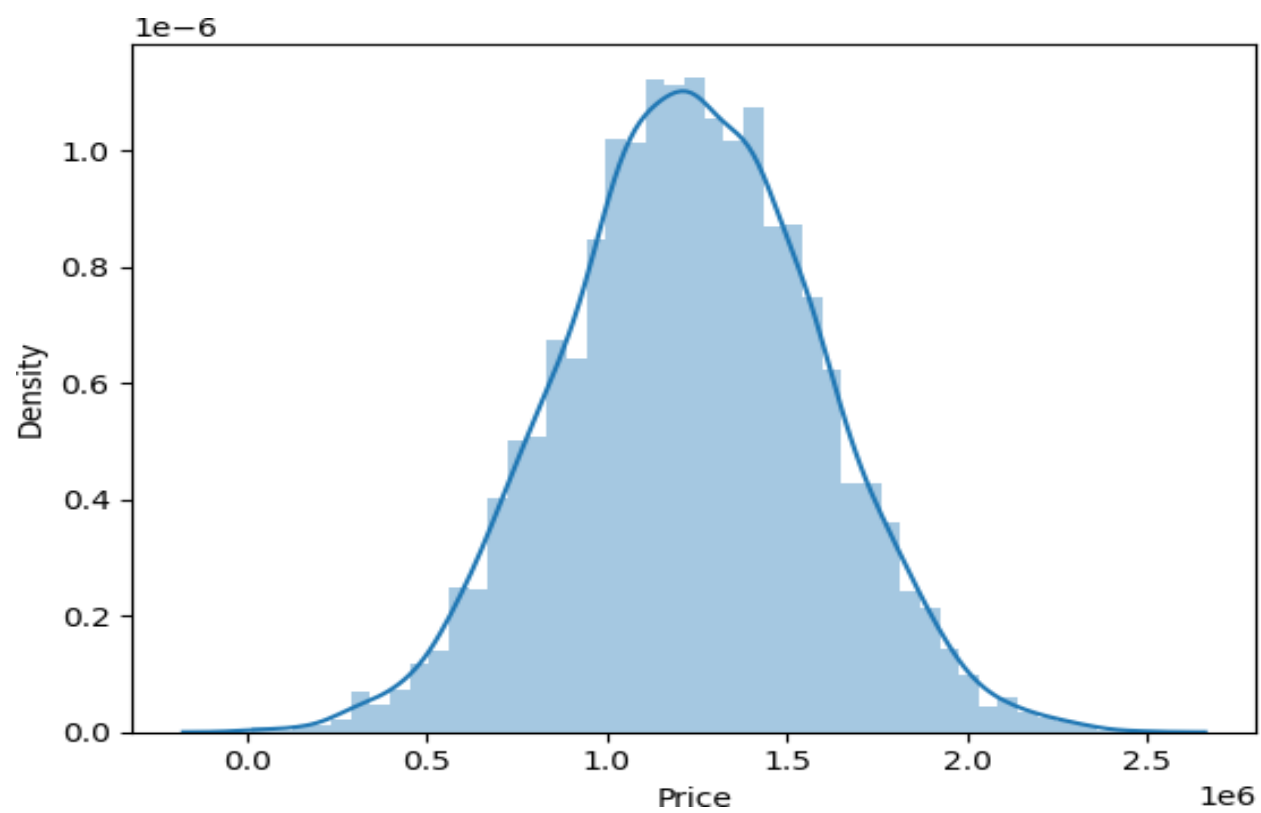


Here we plotted a Scatterplot for average area income.



Average Area Income

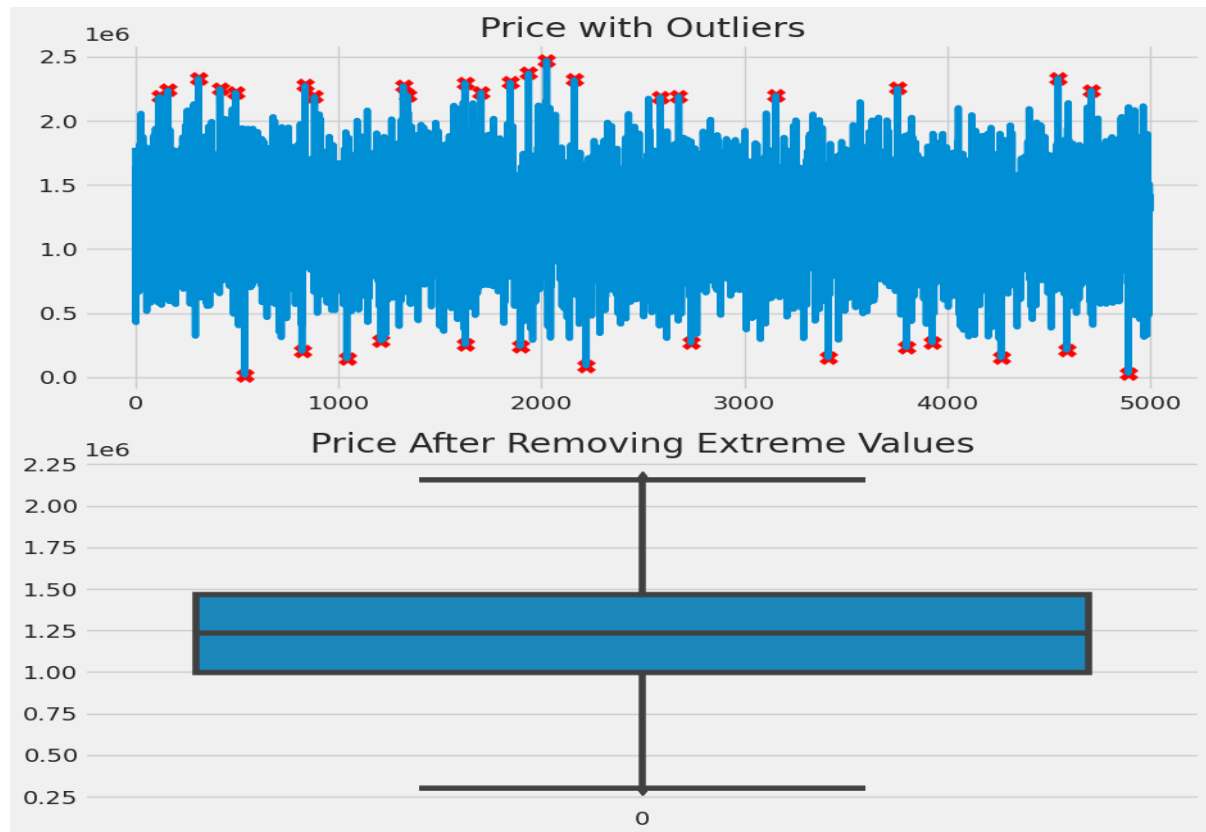
Distant plot for Price and Density in our dataset.



Outlier Analysis

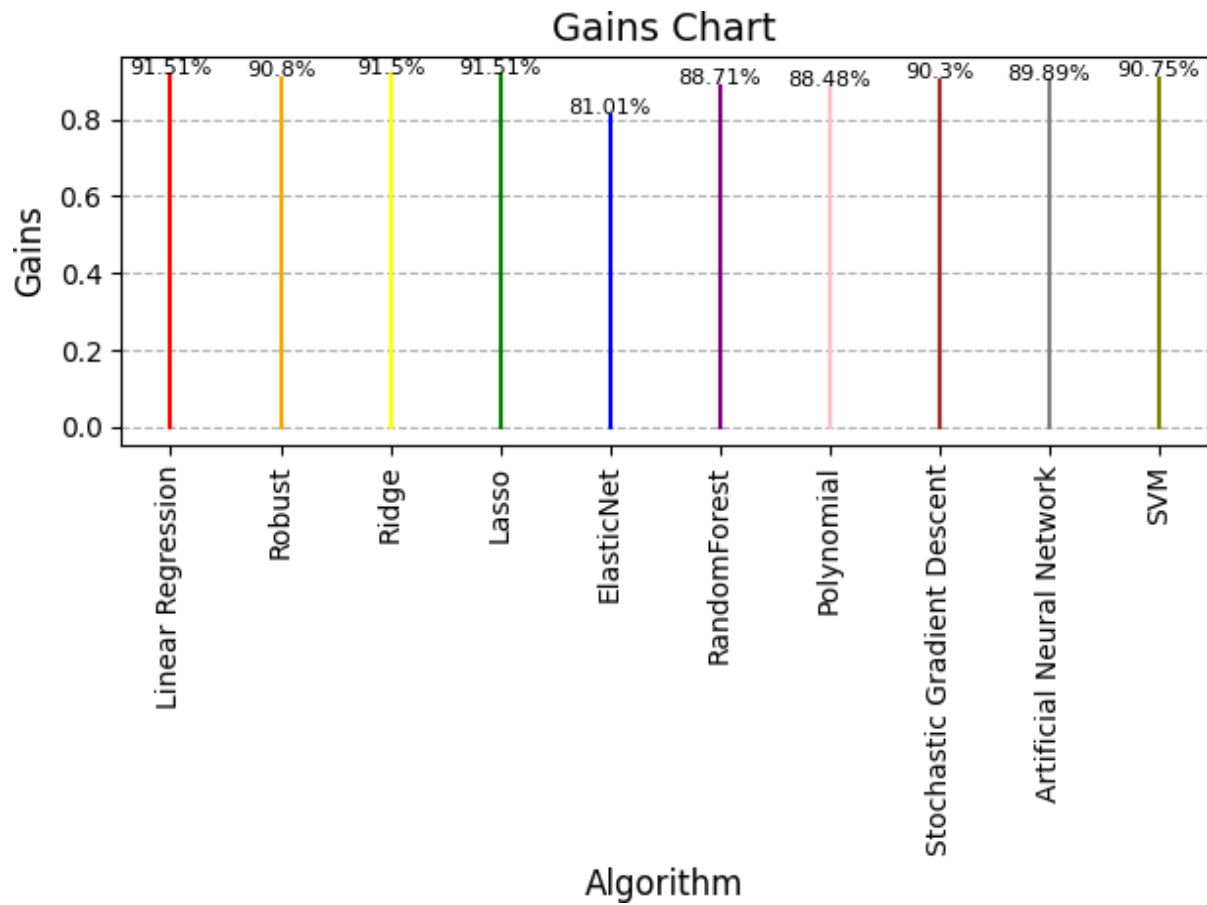
Below we have performed Outlier Analysis in our dataset based on different columns.





Gains Chart

Below we have a gains chart with dotted lines to see the accuracy of all the models that we performed our analysis on.

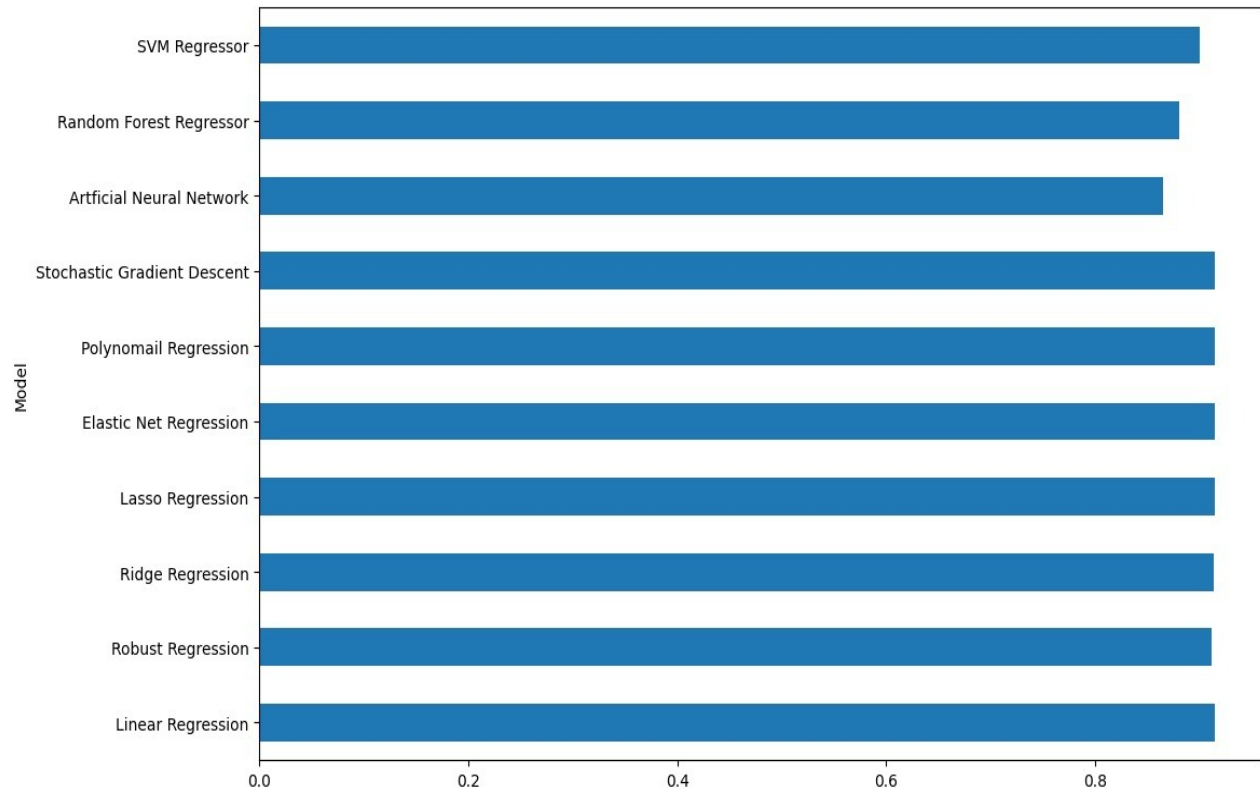


Results

	MAE	MSE	RMSE	R2 Square	Cross Validation
Model					
Linear Regression	81135.566093	1.006842e+10	100341.529545	0.914682	0.917379
Robust Regression	84841.957533	1.109778e+10	105345.989279	0.905959	0.911262
Ridge Regression	81428.648355	1.015327e+10	100763.435337	0.913963	0.917379
Lasso Regression	81135.698517	1.006845e+10	100341.683215	0.914682	0.917379
Elastic Net Regression	81184.431473	1.007805e+10	100389.492321	0.914600	0.879545
Polynomail Regression	81174.518441	1.008198e+10	100409.083243	0.914567	0.000000
Stochastic Gradient Descent	81135.565670	1.006842e+10	100341.528834	0.914682	0.000000
Artificial Neural Network	599600.081205	4.920804e+11	701484.430402	-3.169808	0.000000
Random Forest	94365.592485	1.419648e+10	119148.981672	0.879701	0.000000
SVM Regressor	87205.730510	1.172093e+10	108263.256765	0.900679	0.000000



Models Comparison



Project Outcomes & Recommendations

Project Outcomes

In our project, numerous machine learning algorithms were evaluated for their effectiveness in carrying out a binary classification task. The dataset used for training was normalized, and it was determined that the Linear Regression achieved the best results.

Recommendations

1. To enhance the AI-based system's analytical abilities, natural language processing (NLP) capabilities can be incorporated. At present, the system mainly examines structured data like housing data, but NLP could help analyze unstructured data like customer reviews, leading to the identification of better houses at recommended prices.
2. By utilizing blockchain technology, a secure and transparent system can be established to store and exchange housing data, making real-time feedback implementation feasible.
3. The system can be expanded to other housing search applications by scaling it up.

Appendix: Python Code for use case study

```
Ipip ioetall -g dvgVot
```

- Import Libraries

```
'impom matplotlib.pyplot ae pIt

miatpJotlib iodine

f eve.eet_styIe('whitegrid')
f gIi.etyie:use('fivetAirtyeight")}
```

• Check out the Data

```
USAboueing - pd.read_csv('/content/u°A_Bousing.ceo'
'USAboueing.bead()
```

	ny . area In rose	Avg . Rrea Bouse Sge	Dwg. sea or oz Bourns	Avg . area NuaBer o* eeacz cms	Azea Population
0	79545.458574	5.682861	7.0061 BB	4.09	23086.800503 1.059034e+06
1	78248.642455	6.002900	0.730BZT	3.OB	40a 7 T 072174 T . ñ0b8B1 e +OO
2	012B7.0071 70	5.B00B00	B.5127Z7	?':13	308BZ.1 b9400 T.0@B9 +DB
3	63345.240046	7.1B8356	5.586729	3.26	34310.242831 1.260617e+06
4	?0BBZ. T87ZZ0	5.040555	7.B393B8	4.23	26354.109472 6.309435e+05

```
USAboueing.info()
```

```
<cTae 'paodae.core.tram.DataFrame'>
Baogelndew: a0ns entries, 0 to 4999
Dara eoluene (:otaI - coluoms>:
#   rolumm                               Non-Null Count  Dtype
0   Avg. Lea•income                      a0as oon-ouJl    float64
i   Avg. sea Bouse Age                   s0a0 not-ouz1    float64
?   Avg. •Veabumber• of Boome           Saad noo-oull    float64
z   Avg. sea Pmsber of Beéroors         asaa now-ouIl    float64
   Area Po iation                      aa0s oon-ou2l    Iboa-6
a   Price                              a0aa now-ouJl    float64
6   Addreee                            aa0s oon-ou2l    ob$eM
dtypes: float64(6), object(1)
memory usage: 273.6+ KB
```

```
USWoueing.describe }
```


	Avg. Area Income	Avg. Area Bouse Age	Avg. Area Bumber of Rooms	iv@. area Nuaber o* Bedrooms	Area Population	
count	5000.0JX+0&t0	5000.000000	5000.000000	T00 0D0000	5000.000000	5.000000e+03
mean	f1B5BZI08BB#	5.977222	B.88T?B2	3.OBE Bg0	30T 63.51B05B	1.2Z3073a+OB
std	100b7.9B1214	0.981456	1.005833	1.Z341S7	9BB5.850114	3.fi31178e+0b
min	1T7B6.0g t1 BO	2.644304	3W01B4	2 D00000	1 7Z B1O0BB	1.fiB38B8 e+04
25%	B1480.5023B8	5.3ZZ283	6.299250	3.160000	39403.8Z8702	8.B7?771e+05

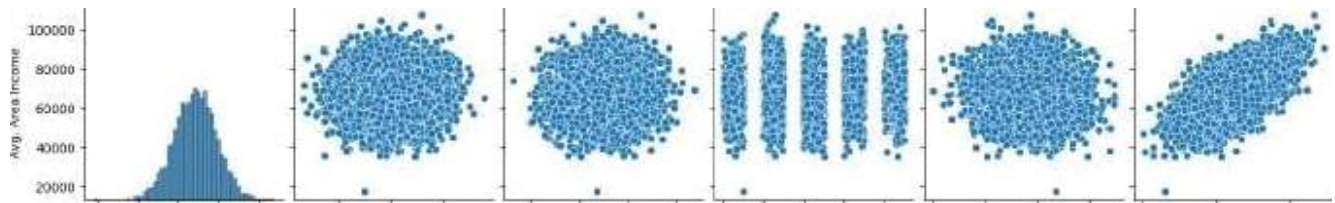
USAboueing.coIrmoe

```
Index({'Avg. Area Income', 'Avg. Area Rouse Age', 'Avg. Area Muober of Roome',  
      'avg. sea number of Sedronme', 'Area Pogulatiao', 'Price', 'Address'],  
      dtype='object')
```

Exploratory Data Analysis (EOA)

sns.pairplot (vSAhousing)

eea5orn: egri0.PqifsriD:at 0p7IaCJéf##YIDi



ipi'g ioetall éeplot

```
Looking in iodeieg: https://gypi.mg/eimply, Atips: {/us-python pig dev/colab-wbeeJs/pubJic/eimplys
Requirement already aatiefied: •bvplét in /Car/local{lr/ honJ.v/died-pacxagee }a. #:s)
Bequiremeot already aatiefied: •pac ging io /uer/locaV/lib/pythnoos.9/dist-pac gee jIrom.hsplot) (z3:i)
maquiremeot already zatiefied: ouxpy>-i.ñ'io /uer{lqcdl/lib/python3.9/dist-pactagee (fzom ñwplotj (1.:z.4)
Bsquiremeot already zntiefied. colorcet -z io /uer/local/J /pytbénz. /diât-paciagea (IN hvplot) (3.a.i)
stir<ml amea »aci^ri>a:toxic•,i.ii.a'in/ r'/iooai/is0/Pmu>na.s/ai^t-Pan 4e<r,oma root fi.is.a)
Bsquiremeot already zntiefied: eJ>-a:i.D io /usr•/IocaJ/iib{g/âond.4{diei-pâctages (from bvplot) #4.I•.A#
aquiremeot already satisfied: •boxeh>-i.A.D'io /uer/local/lib/g#bona.&/diet-pactagee (from hvplot) (?..##
ant a eaa,ea<i^zied: a, in /o'r/local;2u>/ taoni.e/ai t-Mac ge, lrrora rot <i.s.zi
Bequiremeot alr'eady aatiefied: •parâm>ii.s:D'io /ñer{Iocal/lib{g/bood.9/dizt-paciagee•(from hvplot) (?..IA.D#
maq iremeot already zatiefiedz.gilt-7.i.D io /uar/Iocal/lib{p/bond.s{diri-paqiager (from boke>-i.a.é- bvplot
miquirement already aatiefied: •A io9exfex s9-2.ID.A in' iusr/iwa17lib/g bond.4/éiet-gac aqñe (from boreé>-i'
Bequiremeot alrgady eatiefied: tornado>-s.i in /uer/local/lib{p/boo3. {dirt=pactager (from bore>-1.a.0>-bvplot)
Requirement already aatiefied: •PyyI>•J.ié io /uer{léchl/lib/p/bomd.9/dizt-pactagea (from boreb-1.0.#>-beplot]
S irm<m'» am>ad, a<i,ri>a:•?kojaz>-z.'io>^r/inoai/iib/sy<mo3.:7di^ear aga^(ir<m trem-i.c'.a-• im')
Requirement already aatiefied: •gyin>-a.4 in /uer/Voñal/lib{gyt6oo3:9/dist-pactagea #Irom coloicét>-z>6vpldzj
Bequiremeot already aatiefied: •gywis-oomom>-o:7.4'in' /uei/local/lib/ bonsi9/éiei-packages from bolie me>-i:u
maq iremeot already zatiefiedz.g hpn-datertin>-4.8:i io /uer{local/lib7pythons.9/di'et-pactagea (fréepandas- bop:
Bsquiremeot alreddy zntiefied. p/ -4o20.i io /uer{Iécal/lib{p/hdo .S{dizt-pattagez (fret paiflâs- hvplot) (?022'
B#quiremeot alrgady eatiefied: retuptoois>-4z in /usr/Iécal/lib/g/bon3.9/died-gacraquea (from paneU-a,iI.D- dvpi)
Bsquiremeot already zntiefied: car 0wn in /uer/Vocal/lib/actions.g/diet-packagee {from panel>-D.ii.0ññvpIo't) a'
Bequiremeot already aatiefied: •bVeañh in /usr/i al{Iib/pythooJ.v/diet-packages ;from'paoel>-0:ii.0- hvglot) (s.a
B#quieiment alieadyeatiefied: t -d,4#0 io' /usr/2ocal/iib{p/hood.Y/dizt-pactager (from pñel>-0.ii.s>-beplot
Bequiremeot alr'eady aatiefied: •iequeete in /uer/local/lib/pythons.9/died pac 'gee {from'panel>-D.1l.0>bvpioA). <T
mo< ameaar>a<i^éiea:•m xueb%ia>-...o i^/v /Aai/nb/g bona. /di,i-pack'<m, (rr<mnñjat>-2..•meb:
s irm<mo< am>aay »a<i,zied..>i -i.s in/u,rJl'ooai/in/atons.g/ai^c-Pao 4», <*wn hon<eu*il<-z.'i-apr
B#quiremeot alrgady eatiefied: beocodinge:in /uer/locaT/lib{p/bood. {dirt=pactager:(from bleach>-panel.e.li.0•
Requirement already aatiefied: •ieportlib-ñétadata>-4.4 iâ /usr{loc#1/lib{t/éooa• /diat-pactager (from marxxdvo->ç
aquiremdot Already saCiefied: •vzVlibJ<I.2?>,-i.c1.i io /uer/local/lib/pythbn3.9/dist-pactagee (from margoertzn.pa
Bñqu nf aliéady »atieflea: •châreet-nor'oaliser•i.n.0 io /o'ar/local/lib/g bo'na.s/diet-pacraeger (mhm requeeze
miquirement already aatiefied: •idoa 4>>-z.a •io /usr{Iocal/lib/g/honu.9/diét-pactager from regreeie-panel-0.ii'
maquiremeot already zatiefiedz cemifi>-2ei7. „i? in '/rer/locaV/lib/pytAop3:é/diet-p kagee lfrom.requeztz-paoel:
S m<mo< am>aay »a<i,ziea..i -c's in /u,r/iooai/z '7e,tbno,.s70iac-gas», f**nm me>rciit-wtaaa>-'.',--
```

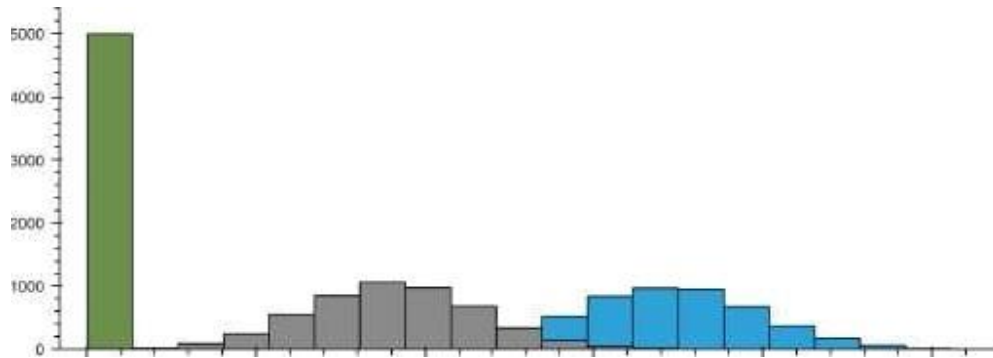
pig ingtall boLoei'ggs

```
Looking iz.indexesz btt // i ,htpiz://ua-python.pkgrdew/Colab-wbeeIs/public?single/
ñemeo< aUeadraa<^flea: boloei ié /uer/local/iib/pythboa.9 digt-paaxagee {l.is,»»)
Bsquiremeot alreddy zntiefied: jsckâging i'o /uer{Iocal/lib/pythboos.9/disi- agea {from hoiovieee) pi3.i}
o ire<mo ahead aaci^fied: oomcet in./u,r/looiailn/g,t one.s./dist-outage,,*roml ivieea) (z.o. }
x irm<mo< aueaay »a<i,ziea. psi>-com,m<-o.z.a in' /»,r/Iocal/iim/ t>o>»v/ai^<-elkage< <re roiwi a> (z.z
Bequiremeot already aatiefied: •paoel>-&.i3,i io /uer/local/lib/pyibon3:»/died-pactagee (from boloviewe) {s:i4:. )
B#quieiment alieadyeatiefied: é -i.0 io /uer/lo'caI/l /pytbooz.9/dist-#c2agea (Irtm<holoviewe) li.z?..4)
Bequiremeot alr'eady aatiefied: •param>?e, >-i.9:3 in' /uer/locaV/Iib/pythons.v/diet-gackagee {froe hoJovise) (s:i3
maquiremeot already zatiefied: •pandae>-4.é.:0 in boar/local/lib/ bon3>s{diet-gackagee {from holovi ) {i.s:3
Bsquiremeot alreddy zntiefied: p/bon-dâteñtil>-2.8.i io /ier/iécal/iib/pytbpn3%>s/di'et-phciâges (fmm pandat>-0.
aquiremeot already satisfied: p/•s>-?&i0.i 'io /uer/local/.iib/g#bon3.&/diet-pactagee (from pandas>-0?>s.s>-boioiv.
Requirement already aatiefied: •oar ñwn io' /uer/local/lib/pytñgza.9/diet pac%gee {froe panels>-D.i3.i>-boloviepe)
aquiremdot Already saCiefied: •bIeach in {uar/Jocai/lib/pjthof3.s/died-pactagee \froe gsie'l'i0.is.i>-ho'lovieue) '
nf añéad, »a<i^lea: 'aéiupCoois>- z in./osr/I ai/iib/ bon3:9/digt-paenaagee (from'pael>a:ii.i>->doio
miquirement aVr'eady aatiefied: •tqda>-a:•#„D. •io /usr{Iocal/lib7g/hon•3.9{diét-pactager" from panel>-0.i3.i>-hñlovi
maquiremeot already zatiefied: debt?..sé, >-Y.a.0''ig /uer7locaI/Vib/pptbons.9/diet-pactagee {from panel>-D.i3:i>-
Bsquiremeot alreddy zntiefied: typlog-extenside in /zrr/local{lib/pythons.s/died-pacraeger jfroe gaoel>-0.la.i>-h
B#quiremeot alrgady eatiefied: .requegte'in' /uer/local/lib/pytbogs.9/diet-pacx°4°° f*1% panei9>-D.13.i>-Foioviewe)
Requirement already aatiefied: •gyin>-0.4 io /uer/Voñal/lib{gyt6oo3:9/dist-pactagea #Irom pangl>•0.i3:i>->âoIoei%
Bequiremeot already aatiefied: •P I>->.iD'io /ñer{Iocal/lib{g/boo3.9/dizt-paciagee•(from boxeb<?..sr0, >-i.4:~0->ç
B#quieiment alieadyeatiefied: t ado>-3.i io /usr{IocaJ/lib{g/hood.9/dizt-pactagee (from bdbreiz.S.0, >-i.a.0w
miquirement already aatiefied: •jiVlom>-7.i.D io /isr/Ihcal/lib/pytbona.s/daai-pac/ager (frñm boxeb<?..a.s> >-...0-
maquiremeot already zatiefied: J}oja>->'.s'io /ugr{lqcal/lib/python3.9/dist-pactagee (fzom boreb J.a.o> >-2.4.4>->pt
o irm<mo< ameady aa<i visa. air-i.c in /ur/l'ocai i / tñonz.g/area-sac ge <*wn hon<euil<-z.'a.i-ver
aquiremeot already satisfied: •webeocodinga'io /uer/local/.lib/g#bon3.&/diet-pactagee (from ble b>-panel>-D.13.J•
Requirement already aatiefied: •ieportlib-ñátadata>-4.4 iâ /usr{local/lib{p/boo3. /dizt-pactager (from mar vo->ç
Bequiremeot already aatiefied: •rr1lib3<i.i?> >-i.i.i io /uer/local/Iib/gytbooz.9/dist-gackagee (from'requeetañpa
```

requirement already satisfied: idna<u,>-i.s in ?uer/locaJ/lib g/bond.9 diet-pacaagee (from regreeie->ganeI*-0.id

requirement already satisfied: cbareet-noreaIizer--°.n.n in /usr/local/lib/p)vdona. /die:-gacrages (from requeteet-
 Bsquirement already eataefied: ceñifi>-?0i?..17 in /rer/local/lib/gytZons.9/diet-pacxagee bros requests- pael:
 requirement already satisfied: sippo-a.S iz /uer/loc*I/Jib/pytboos.9ldist-pacxagee <from iigortlib-metadata=-#.:-:
 Bsquirement already eataefied: :mrkug&aIe>-z.0 in /uer/Iocal/lib/p hone.»/died-pacagee (from Jizja?>-?. -oleh*

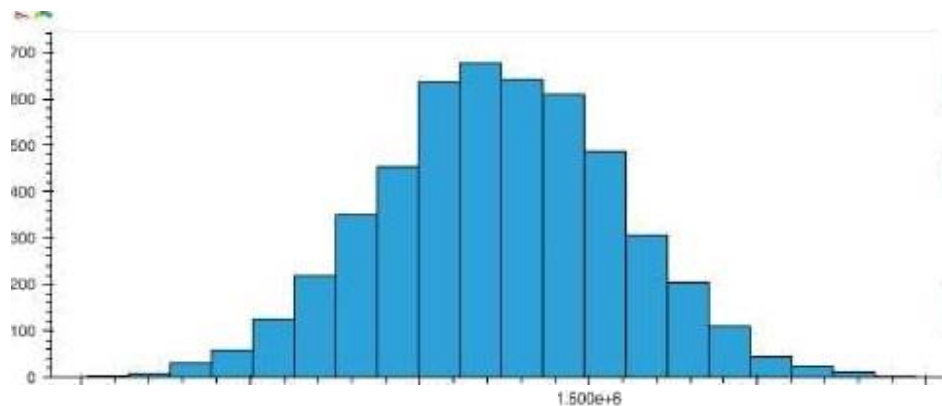
```
import holoviewee ae hv
hv.emeneion('bekeh')
UsVoueing.hvoIoz.hist{by-'Prire', srbplots-Falee, ridtb-i000)
```



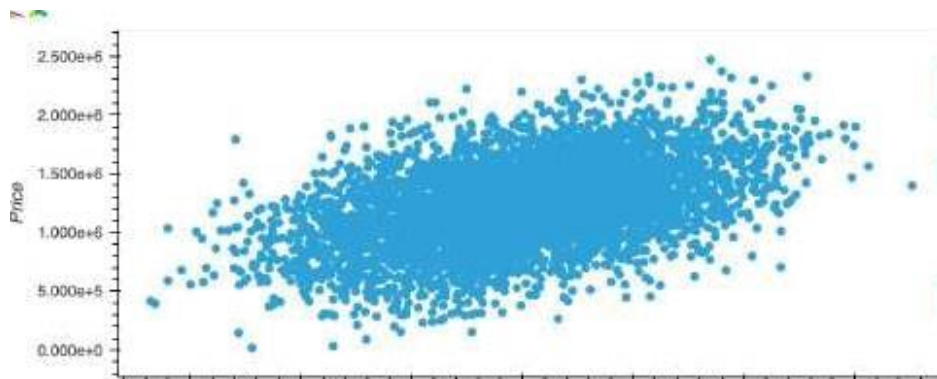
Arg. 4ree in sme

Avg. Area Number r+4 Rooms
 Avg. Area Numbar o4 Bedroom*

```
import holoviewe ae bv
hv.esxension('boxeh')
UsA oueing.bvpIm.hist{"Price"1
```



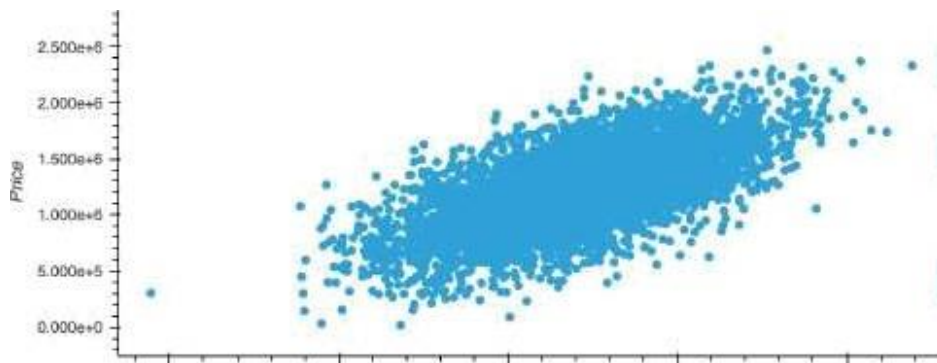
```
imprnr holommn B lv
bv.eaxeneion(' eh')
USAboueing.bvpIm.ecatterIx-'Avg. Area aouse &qe', y-'Price'1
```




```

impnm holomewn n hv
bv.erreneioo(' eñ')
USAAoueing.bvpIm.scatter(x='Avg. Area inrome', y='Price')

```



```

impom holovieeee ae hv
bv.emeneion('bokeh'{
USAboueing.coirmoe

```

```

Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Dumber of Rorme',
      'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Addresser'],
      dtype='object';

```

```

sts.beatiap(nS.housing.corr(), anoot=True)

```

<âpyt from—input-JS\$—4366Be6a5 o6 £>:i: Tut uroWarni rig: The fie£ as £t va£ue of nunz r:icyonly in Datagram. norr in dspreca-
sue. tleatarap (JJSiif onsing.corr)[, an not-True)



• Logistic Regression Model

X and y arrays

```
z = DsAhousiog({'Avg. sea lncoom', 'Avg. Area moug's Age', 'Avq. Area number oI Boome',
               'Avg. Area Puerder o'f Bedrooms', 'Area Population' })
y = D'Ahousioq{'Price' }
```

• Train Test Split

flow, let's sport the data into a training set and a testing set. We will train our mode4 on the trBiNking set and then use the test set to.eyalual

```
from eklearn.modeJ_eelectioo import irain_teet_eplic

z_train x_tert, y_Craio, y_<eet - train_test_sp?it x, y, teei_size-n.s randoz_eCate-•?#

bros eilearn iaport metrics
froe exlearn.codel_eelectioo import croee_val_rcnre

del croee_val(sedel):
    pred=croee_val_score(model, J, y, cr-i0)
    return•gred.sean{ }

def•pmnt_evalraie true, predicted):
    mae=metrics.mean_abeolute_error(true, predicted)
    mee=met:ice.mean_rquared_error(true, gred e0
    imse - np.rqr: metrics.mean_sgrareñ_error(true gre cued))
    ri_equare - :metrics.r•_ecoreftrue, predicted)
    print 'I: ', mee
    print 'msx:', mee
    print('RMSE:', rmse)
    pre nt 'R*squar e', • z square [
    print ' _____' )

def•maluateltrue, predictedlr
    mae - met:ice.meao_abrolute_er•rgr(true, predicted)
    mee=metriJe v» rquar#d_er'ror(true, gredi:ned
    imse=np.rqr:metrics.mean_sgrareñ_error(true gre cued))
    ri_equare -:metrics.r•_ecoreftrue, predicted)
    return mae, mee, rmse, ri_equare
```

• Preparing Data For Linear Regression

- Linger As8umpt .
- Oeussien Oistzdzutiora.

```
from eklearn.preprocessioq import Staoardscaler
froe exlearn.gipeJine izpomPigelize

gipelioe - Pipeline({
    ('etf_ecalar', stand dscaler{ })

¥_traio - pipeline.fit_traoeform x_train]
x_teei - pipeline.transform(a_test)
```

• Linear Regression

```
from sklearn.linear_model import LinearRegression
```

```
lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)
```

```
• LinearRegression()
  LinearRegression()
```

• Model Evaluation

Let's evaluate the model by checking our coefficients and how we can interpret them.

```
lin_reg.intercept_
```

```
1228219.1492415662
```

```
coef_df = pd.DataFrame(lin_reg.coef_, x.columns, columns=['Coefficient'])
```

	Coefficient
Avg. Area Income	232d7B7248A3
Avg. Area House Age	163B41.046583
Avg. Area Number of Rooms	121110.5s547B
Avg. Area Number of Bedrooms	1289Z.8INt1B
Area Population	1f1ZZ234?37I

Double-click (or enter) to edit

• Predictions from our Model

Let's grab predictions off our test set and see how well it did!

```
pred = lin_reg.predict(X_test)
```

```
pd.DataFrame({'True Value': y_test, 'Predicted Value': pred}).head()
```

Residual Histogram

```
pd.DataFrame({'Residual': y_test - pred}).head()
```

• Regression Evaluation Metrics

Here are three common evaluation metrics for regression problems:

- Mean Absolute Error (MAE) is the mean of the absolute value of the errors.
- Mean Squared Error (MSE) is the mean of the squared errors.

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Comparing these metrics:

- MAE is the easiest to understand, because it's the average error.
- MSE is more popular than MAE because MSE 'punishes' larger errors, which tends to be useful in the real world.
- RMSE is even more popular than MSE. RMSE is interpretable in the 'y' units.

API of these are Yes Nelsons, because we want to minimize them.

```
test = model.predict(a_test)
train = model.predict(a_train)
```

```
prior('Test set evaluation:\n_____')
print('Train set evaluation:\n_____')
prior_evaluation(train, test)
```

```
result = pd.DataFrame(data={ 'Train Regression', *evaluation(y_test, test), cross_validation()}).
    columns=['Train', 'Test', 'Cross Validation']
```

Test set evaluation:

```
MAE: 81135.56609336878
MSE: 10068422551.40088
RMSE: 100341.52954485436
R2 Score 0.9146818498754016
```

Train set evaluation:

```
MAE: 81480.49972774892
MSE: 10287043161.197224
RMSE: 101425.06180031257
R2 Score 0.91929667975527
```

• Robust Regression

• Random Sample Consensus - RANSAC

```
from sklearn.linear_model import RANSACRegressor

model = RANSACRegressor(backend='linear', max_trials=100)
model.fit(x_train, y_train)

test = model.predict(a_test)
train = model.predict(a_train)

prior('Test set evaluation:\n_____')
prior_evaluation(train, test)
print('-----')
```

```
prior('Train rei evaluation:\n _____')
prioievaloate(y_irain, train#red)

rerrlie_df_z - pd.DafaFrame(data-{"mnbuei egreeeion", *evaluate(y_teet, test_reds, croee_val{RABdACRegreseor()1}},
                                colusme-{'xodel', 'xA&', 'mss', 'VSs', 'B? square', "Cross validation'})
reeulie_df - reeuJte_df.append(resrIte_df_z, ignore_index=True)

/uer/local/lib/pytñoos.9/diet-oacxagea7srlearn/linear_eodel?_razsac.gy:#4?: ruiurewarning: 'base_eeimator' wae r'
warnirgs.who
Ceei eet evaluation:

_____

MSE: 10829564871.937178
RMSE: 104065.195295724
R2 Square 0.9082320555368967

_____

Crain eet evaluation:

_____
: 83910.60259G69 117
MSE:
RMSE: 104552.13670255442
R2 Square 0.9142456773615641

<ipythoo-input-i i-zoraacbdzacc:i : Futurewarning: Cbs 'frame.append method ie deprecated and will be removed fri
rea<lts_df - reulte_df.appendlresults_df_?, ignore_index=True)
```

• Ridge Regression

```
from eklearn.Iioear_mddeI import Ridge

model - Ridge alpa-ieo, solver-'cdoleety!, too-0.a00i, random_state- c)
model.fit x_traio, y_traiz)
pred - moéel.predict{«_teei)

teet ed-modea.predict(a_ieet>
'traio#red -moéeJ.predict(a_zrain)

prior('Teet set.evaluation:\n _____')
.grio:_evaluate(y_ieet, tee:bred)
prior('_____')
prior('Train rei evaluation:\n _____')
prioievaloate(y_irain, train#red)

rerrlie_df_z - pd.DafaFrame(data-{"Bidge Regreeeeion", *evaluatef¥_""; rest i*d1, cr e_oal{Ridge()}{},
                                colusme-{'xodel', 'xA&', 'mss', 'VSs', 'B? square', "Cross validation'})
reeulie_df - reeuJte_df.append(resrIte_df_z, ignore_index=True)

Cee- eet evaluation:

_____
: 8142B.64B3553*+°136.
ME : JOL5'326990d.B9°6D9
CASE: I OOd63.43538689494
820quare D.9139G286VAA64G09

_____

main eet evaJuatioo:

_____
IU+E:819.7E.39d58585fi' DT

V::E: iolsvs.46l s:i9:s?
a? souare n.9izaa44F?::#?444

<iphon-iopui-i 4-dDaAivsef:7 >'zt: xtturewarning: .ibd*frame.append method ie degrecated and will be removed frx
results_df - results_df.append(results_df_2, ignore_index=True)
```

• LASSO Regression

4/23/23, 9:26 PM

HoudoçR:wrRrdkâea-Hamid-Cobberr

```
from sklearn.linear_model import LogisticRegression
```

```
model = LogisticRegression(
    random_state=0,
    solver='lbfgs',
    multi_class='ovr',
    tol=1e-5,
    verbose=1,
    warm_start=True,
    selection='random',
    random_state=0)
model.fit(x_train, y_train)
```

```
y_test_pred = model.predict(x_test)
y_train_pred = model.predict(x_train)
```

```
print('Test set evaluation:\n')
print('Train set evaluation:\n')
print('Cross-validation score:\n')
print('Coefficient of determination (R^2):')
print('Adjusted R^2:')
print('F-score:')
print('Area under the ROC curve:')
print('Log loss:')
print('Brier score:')
print('Mean absolute error (MAE):')
```

```
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, \
    roc_auc_score, brier_score_loss, mean_absolute_error, mean_squared_error, \
    log_loss, r2_score, adjusted_r2_score
```

Test set evaluation:

```
MAE: 0.118443147330944
MSE: 0.060546847014997
RMSE: 0.246080944014997
R2: 0.9146002670381437
```

Train set evaluation:

```
MAE: 0.08147330944
MSE: 0.02813196614997
RMSE: 0.16771814997
R2: 0.992220997
```

```
<ipython-input-14s-abtb:aadidi?>: FutureWarning: Cbs frame.append method is deprecated and will be removed from
result_df = results_df.append(results_df, ignore_index=True)
```

Elastic Net

```
from sklearn.linear_model import ElasticNet
```

```
model = ElasticNet(alpha=0.0001, l1_ratio=0.5, random_state=0)
model.fit(x_train, y_train)
```

```
y_test_pred = model.predict(x_test)
y_train_pred = model.predict(x_train)
```

```
print('Test set evaluation:\n')
print('Train set evaluation:\n')
print('Cross-validation score:\n')
print('Coefficient of determination (R^2):')
print('Adjusted R^2:')
print('F-score:')
print('Area under the ROC curve:')
print('Log loss:')
print('Brier score:')
print('Mean absolute error (MAE):')
```

```
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, \
    roc_auc_score, brier_score_loss, mean_absolute_error, mean_squared_error, \
    log_loss, r2_score, adjusted_r2_score
```

Test set evaluation:

```
MAE: 0.118443147330944
MSE: 0.060546847014997
RMSE: 0.246080944014997
R2: 0.9146002670381437
```

Train set evaluation:

: 8 157 7 . B 8 B 3 J 53 17 BI

RMSE: 101485.34351373828
R2 Square 0.9192027001474953

<iphon-input-1 <-zae?fee?c#i4>'z?: Futurewx*ning• Cue frame.append method is deprecated and will be removed from
results_df - results_df.append(results_df_?, ignore_index=True)

Polynomial Regression

```
from sklearn.preprocessing import PolynomialFeatures

poly_reg = PolynomialFeatures(degree=3)

x_train_poly = poly_reg.fit_transform(x_train)
x_test_poly = poly_reg.transform(x_test)

lin_reg = LinearRegression()
lin_reg.fit(x_train_poly, y_train)

test_results = lin_reg.predict(x_test_poly)
train_results = lin_reg.predict(x_train_poly)

print('Test set evaluation:\n')
print('Train set evaluation:\n')
print('R2 score: %f' % r2_score(y_test, test_results))
print('R2 score: %f' % r2_score(y_train, train_results))

# Evaluate the model
from sklearn.metrics import mean_squared_error, r2_score

test_results = lin_reg.predict(x_test_poly)
train_results = lin_reg.predict(x_train_poly)

mse = mean_squared_error(y_test, test_results)
r2 = r2_score(y_test, test_results)

print('Test set evaluation:\n')
print('Train set evaluation:\n')
print('R2 score: %f' % r2_score(y_test, test_results))
print('R2 score: %f' % r2_score(y_train, train_results))
```

see: test evaluation:

IU+E: 8 1 L 7 I . B J 8 1 4 1 L 9 6 9 S

RMSE: 100409.08324260656
R2 Square 0.914566932419506

Crain test evaluation:

MSE: 101323.67517519198
R2 Square 0.9194599187853729

<ipython-input-1 +-FDIbcieo0 iu-:?: Future warning: CFE frame.append method is deprecated and will be removed from
results_df - results_df.append(results_df_?, ignore_index=True)

• Stochastic Gradient Descent

```
from sklearn.linear_model import SGDRegressor

sgd_reg = SGDRegressor(max_iter=1000, tol=1e-4, verbose=1)
sgd_reg.fit(x_train, y_train)

# Test the model
test_results = sgd_reg.predict(x_test)
train_results = sgd_reg.predict(x_train)

print('Test set evaluation:\n')
print('Train set evaluation:\n')
print('R2 score: %f' % r2_score(y_test, test_results))
print('R2 score: %f' % r2_score(y_train, train_results))
```

```
reruJie_df_z = pd.DataFrame(data={'stochastic Gradient Descent', 'evaluate_y_tee:', test red ,e ;
                                col=I'xodel', 'xA&', 'mrs', 'VSs', 'B° square', 'cross va?idatioo' })
rerrie_df = reeulte_df.agpené resrIte_df_?, ignore_inéex-Tue]
```

Ceei eet evaluation:

NAE: 81L715.561245706E

RMSE: 100341.52846153275

R2 Square 0.9146818517176497

Crain eet evaJuatioo:

: 8J4BD.49B7J82496

10287043161.210154

RMSE: 101425.06180037631

R2 Square 0.9192986579074511

<ipython-input-i #-4• I?99bi?i#>: FutureWarning: Cbs frame.append method is deprecated and will be removed from
 reeltg_dI = resulte_df.append reeltts_df_?, ignore_index=True)

• •/" Artificial Neural Network

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Input, Dense, Activation, Dropout
from tensorflow.keras.optimizers import Adam
```

```
z_train = np.array(x_train)
z_test = np.array(x_test)
y_train = np.array(y_train)
y_test = np.array(y_test)
```

```
model = Sequential()
```

```
model.add(Dense(10, activation='relu'))
model.add(Dense(10, activation='relu'))
model.add(Dense(10, activation='relu'))
```

```
model.add(Dense(6, activation='relu'))
# model.add(Dense(10, activation='relu'))
```

```
model.add(Dense(10, activation='relu'))
* 0.01, add(Dense(10, activation='relu'))
```

```
model.add(Dense(10, activation='relu'))
model.add(Dense(10, activation='relu'))
model.add(Dense(10, activation='relu'))
```

```
model.compile(optimizer=Adam(), loss='mse')
```

```
r = model.fit(x_train, y_train,
              validation_data=(x_test, y_test),
              batch_size=10,
              epochs=10)
```

```
3500/2500 [-----] - 13s 3ms/step - loss: 1635543154688.0000 - val loss: 1657098731520.0000
```

```
3500/3500 [-----] - 12s 3ms/step - loss: 1635543154688.0000 - val loss: 1657098731520.0000
```

```
Epoch 3/10
3500/3500 [-----] - 11s 3ms/step - loss: 1635543154688.0000 - val loss: 1657098731520.0000
```

```
3500/3500 [-----] - 13s 4ms/step - loss: 1623434723328.0000 - val loss: 1637010898944.0000
```

```
3500/3500 [-----] - 13s 4ms/step - loss: 1599173427200.0000 - val loss: 1601585020928.0000
```

```
3500/3500 [-----] - 12s 3ms/step - loss: 1547818500096.0000 - val loss: 1531431092224.0000
```

4/23/23, 9:26 PM

HoudogR:wrRrdkâea-Hamid-Cobberr

```
3500/3500 [-----] - ce 3sm/etep - lose: 1 sisa: zissz8.aaoo - val loes: i o6vs sioszz:ao'
Egon ft 8 / L.O
350 0/ 3500 [-----] - isa 4ms/atep - loss: ::siaaoo:ieez.ooou - val loss: no i i*zB:i:oo'

350 0/ 3500 [-----] - 13s 4ms/step - loss: 1053762191360.0000 - val loss: 931060580352.0000
3500/?500 [-----] - ozs cmm/step - lose: z<zaiaz ssoe.aoao - val_loes: ses scazeiDi.ooaa
```

```
pd.DataFrame({'True Valree': y_test, 'Predicted values': gred}).bvgIoi.scatter(x='True valuee', y='Predicted vaieua'l
```

```
.pd.Dair°rame;z.history}
```

```
II 1.03b077e+12 1. 1T0e+IZ
1 1.035543e+12 1.657099e+12
2 1.632863e+12 1.0fs0fl7e+72
0 T.023435e+12 1 7Ot Ie•12
4 I:50B173e+12 1.601585e+12
'i t.5478 T0e+1 Z 1.531431e+12
B 1.6dl.BB3e+12 1.408755e+1 Z
7 1 14 B0e+12 1.3001'4 te+IZ
8 1.053762e+12 9.310808e+11
8 7.4B 7B3e+11 5.957625e+11
```

```
gd.Dat 'rame;i.bietñry}.b Jot.line(y-{'loss', 'val_lose' )
```

```
test ed-model.predict(s_test)
traio#red -:endeJ.predict(a_train)
```

```
prior('Teet aet evaluation:\n_____')
prio_i_evaloate(y_ieet, teai#red)
```

```
grint('Train rei evaluation: n_____')
prio_i_evalzate(y_irain, train#red)
```

```
resulte_df_z - pd.DafaFr #data•({'artfici'aI 8enal metvorx', '*evalraie y_test, teetered), a{ ,
.col.m.m ;'nodél', 'xAc', . 'nss', 'msss', 'ai square', 'croas validation'})
reruJie_df - reeulte_df.appende resulte_df_:, rgoore_index-True
```

```
47/47 [-----] - 0s 2ms/step
mo/iaa ; [-----] - 0s 2ms/step
reef eet evaluations
```

M: "24S96.46009GO566

771856.5495098421
R2 Square -4.048393327894736

Cram eet evaluation:

MSE: 574624782796.526
758040.0931326298
R2 Square -3.5079028487194366

```
<ipython-inguz-id+-bzcbase?ii>:ii: Future arning: CFfe frame.append method ie dégreced and will be removed fri
results_df - resulte_df.appende reeults_df_?, ignore_index-True)
```

• Random Forest

```

from sklearn.ensemble import RandomForestRegressor

rf_reg = RandomForestRegressor(n_estimators=100)
rf_reg.fit(x_train, y_train)

# Test set evaluation
y_pred = rf_reg.predict(x_test)
mse = mean_squared_error(y_test, y_pred)

# Train set evaluation
y_train_pred = rf_reg.predict(x_train)
mse_train = mean_squared_error(y_train, y_train_pred)

# Cross-validation
cv_mse = cross_val_score(rf_reg, x_train, y_train, cv=5)

# Print results
print('Test set evaluation:')
print('MSE: %.4f' % mse)
print('R2: %.4f' % r2_score(y_test, y_pred))

print('Train set evaluation:')
print('MSE: %.4f' % mse_train)
print('R2: %.4f' % r2_score(y_train, y_train_pred))

print('Cross-validation:')
print('Mean MSE: %.4f' % cv_mse)

```

FutureWarning: The following attributes of the RandomForestRegressor class are deprecated and will be removed in version 1.2: feature_importances_, oob_score_, oob_prediction_. Please use the corresponding attributes of the RandomForestClassifier class instead.

- Support Vector Machine

```

from sklearn.svm import SVC

svm_reg = SVC(kernel='rbf', gamma=0.001, C=1.0)
svm_reg.fit(x_train, y_train)

# Test set evaluation
y_pred = svm_reg.predict(x_test)
mse = mean_squared_error(y_test, y_pred)

# Train set evaluation
y_train_pred = svm_reg.predict(x_train)
mse_train = mean_squared_error(y_train, y_train_pred)

# Cross-validation
cv_mse = cross_val_score(svm_reg, x_train, y_train, cv=5)

# Print results
print('Test set evaluation:')
print('MSE: %.4f' % mse)
print('R2: %.4f' % r2_score(y_test, y_pred))

print('Train set evaluation:')
print('MSE: %.4f' % mse_train)
print('R2: %.4f' % r2_score(y_train, y_train_pred))

print('Cross-validation:')
print('Mean MSE: %.4f' % cv_mse)

```

```
MSE: 9363827731.411339
RMSE: 96766.87310960988
R2 Square 0.9265412370487783
```

```
<ipython-input-170-5d3d7f23418c>:17: FutureWarning: The frame.append method is deprecated and will be removed from
results_df = results_df.append(results_df_2, ignore_index=True)
```

```
results_df
```

	Model	MAB	MSB	RMSE	M square	Cross Validation	*
0	Linear Regression	811.0	880.3	1.000842e+10	100341.529545	9.9140BZ	08bB0f1
1	Random Forest	834.2	9.773	1.083068e+10	104065.1B52B6	0.808232	0.81"379
2	Ridge Regression	81438.0	83ns'Laaso	1.015327e+10	100703.43533a	0.9138g3	0.B173?8
3	Regression	11.0	88n17	1.000845e+1B	100041.8B3Z15	0.8140B2	0B?BZ#5
4	Elastic Net Regression	118-4	4314?3	1.007805e+1B	10958B.4B?3Z1	0.914000	0.000000
5	Polynomial Regression	11.0	844t	1.00818Be+10	10540B.0B3243	0.914567	0.000000
6	Stochastic Gradient Descent	11.0	84246	1.000842e+10	100341.5Z84g2	0.914082	0u000000
7	Decision Tree	8kc1.0	000212	1.41748Be+10	118058.ZB5288	0.8?88B#	0.0000D0
8	SVFJ Regression	872'05.7	30010	1.1T2093e+10	10B203.256785	0.900679	0.917379

Models Comparison

```
resrie_df.set_index('Model', inplace=True)
resrie_df['R2 Square'].plot(kind='bar', title='R2 Square')
```

```
<Axes: ylabel='Model':
```



```
import numpy as np
import matplotlib.pyplot as plt
```

```
#define values and corresponding algorithm names
```

```

:va triea - [a.s¥saaayazs4ss7a ,a.9o:ssPa aJa as e z s z.,°a.é t4 se 7 s y s s's sais TT,a.s Jso e oa s 4 a z a'Jls *,
0.9T5BBBB'1%3LV24969,0.8LDD944.7737L64',0.B87g96.74497J698B,0.88d83'T463B29d6-3'6,
0'IB246668*08787"B95,0.907990644*37L66B?,0.898936*2B55d46069,B.90T459407066L2*,
0.7508569551273301, 0.8693855376039462, 0.9150085163619179]
'algoriihm_namer*-[ 'zogirtic Begreazion', 'Robret?', 'TbeiI men', '&idge',• 'zaseo', 'sVaztic8m:': 'sandocForert',
'x6Boosi', 'Adabooei', 'GradieocBoost', 'iigbteSN', 'catBoori', 'Deciaioo free', 'KKN', 'SvM']

f coooert véIree'io percentage ané rouoé'to tvd decimal glacee
vairee t- {round;yal.* i0a, z) for val in valuee{

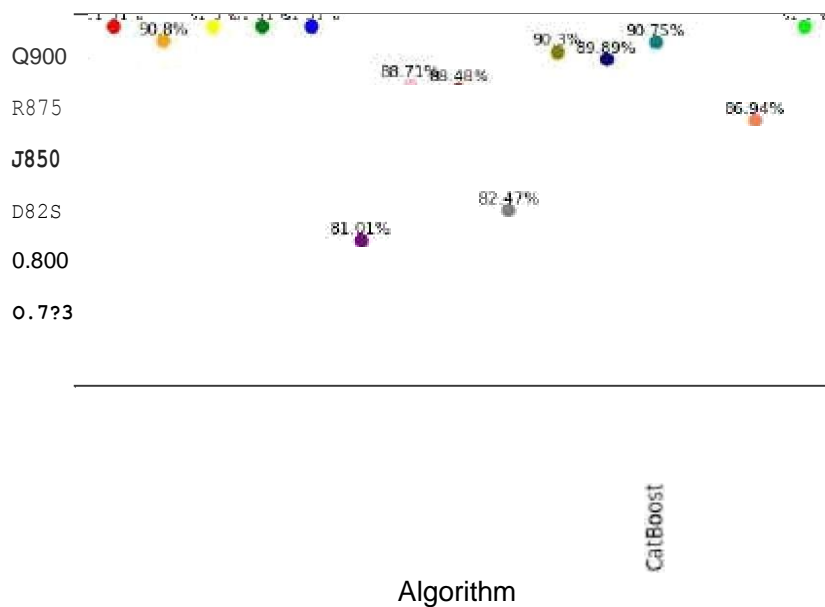
#'creaie a list of colors féieach algorithm
colore' - ['red', 'orange', 'yellom', 'green', 'bJue', '{nrrple' 'gink' 'brovo'; 'gray™, 'oJive',
!navy', 'teal', 'magenta', 'coral', 'lime'

é create the Aft châit•
fig, ax - pit.rubplots()
ior i in range(len valuer)):
    plt.plo't({dlgoritñm_ramee[ij , [valuee[ij], color•coIoie\i} marker-'o'}
    'pJ<.texrlalgorithm_zamee{i{, valree[i]+0:00s :etrçvaIree#ct(i)}+ '% ' ña-'ceñier', footeize-#}
plf.xlabel\ 'AlgoiitM' Ionieiae-I•#
pit.yJabelç' ion', footeise-i?j
pit.:iile 'saioe thart', Ioñteize-i4j
pit.sxicae{rotaiioc-s0>
glt.iight_layoui()

```



Gains Chart



Colab paid products Cancel contracts here

✓ 2s completed at 9:26 PM

