

Project :

Laptop Price Predictor

So, The Main Objective of this project is Predicting The Price of laptops. Using The past data , it Undergoes The Data preprocessing and Exploratory Data analysis . and We Use "Machine Learning" algorithms to predict The model.

The output (dependent variable) is "Continuous data" say what. !
data is in 'Numeric data' , so . we apply Regression models . For prediction.

So, On Observation we can see that if we remove

"GB" From RAM, i can make it as integer value Then after, now same goes with memory, weight for weight i can classify it as floating variable using

The str.replace(). 8GB \Rightarrow 8.

- $df["Ram"] = df["Ram"].str.replace("GB", "")$
Error : can only use str accessor with string values!
- $df["Ram"] = df["Ram"].astype(str).str.replace("GB", "")$
datatype changed
- $df["Weight"] = df["Weight"].astype(str).str.replace("kg", "")$
converting from string = integer for "Ram" column
- $df["Ram"] = df["Ram"].astype("int32")$
Converting from string = float for weight column
- $df["Weight"] = df["Weight"].astype("float32")$

df. head()

Ram	Weight
8	1.37
8	1.34

old value: 8 GB

1.34 Kg.

df.info()

Ram = int32.

Weight = float32

old info = object

old info = object

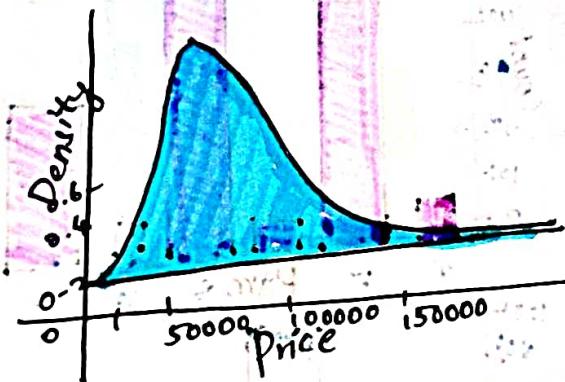
EDA

sns. distplot

(df["Price"])

↳ output variable

dhi distribution plot
edi



Right Skewed.

plotting

Countplots for

"Categorical variables"

* Company

* TypeName

* Ram

* OPSys

Object =

categorical
data

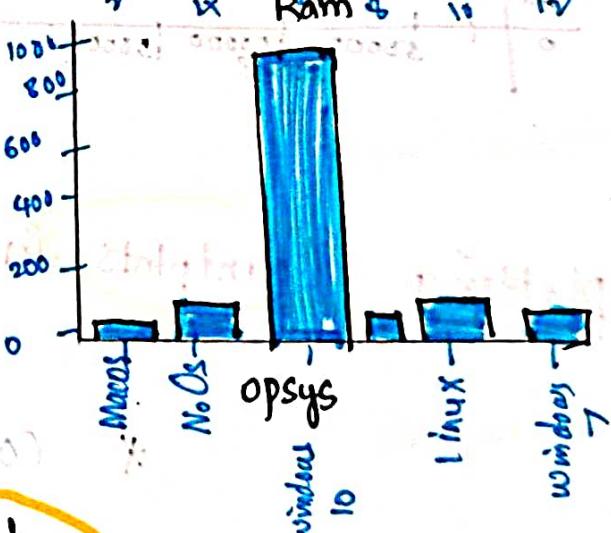
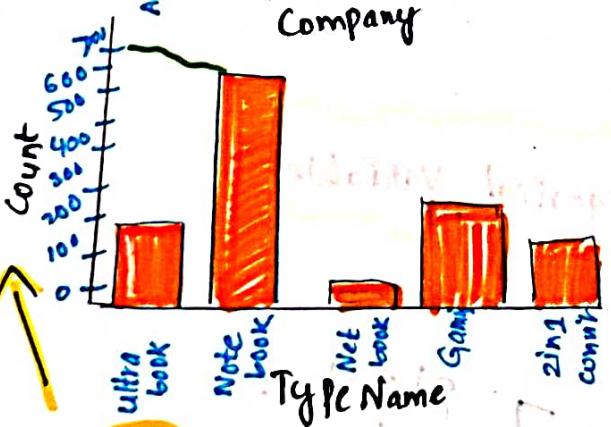
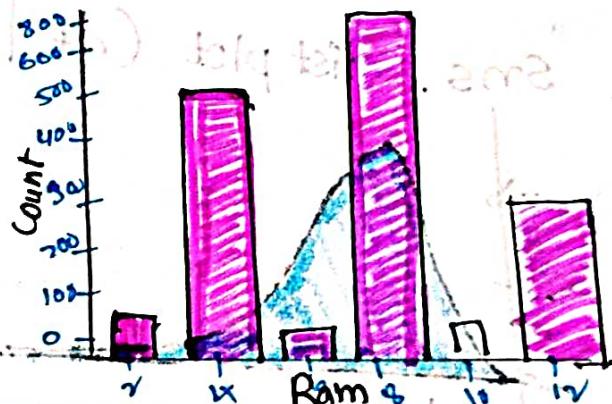
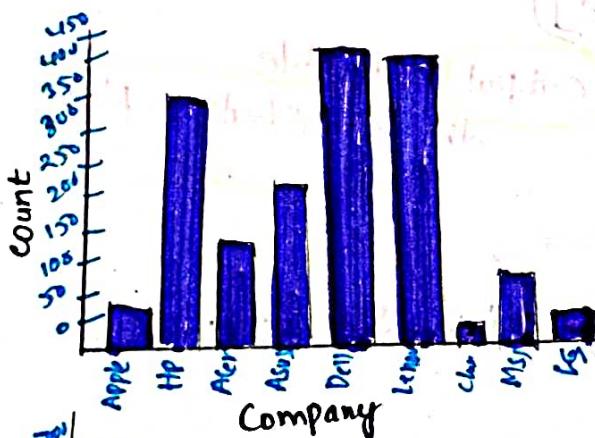
function ← name → column
 def drawplot (col):
 plt. figure (figsize = (15, 7)) \Rightarrow (7, 3)
 only last figure print avuthadi
 sns. Countplot (df [col])
 plt. xticks (rotation = "vertical")
 place vundakku pratt dhaniki, --- (1) 111 \rightarrow Edi easy ga ardam
 to view = ["Company", "TypeName", "Ram", "OpSys"]

for col in to view:

categorical data (object):

drawplot (col)

draw CountPlot.
ani ardham.

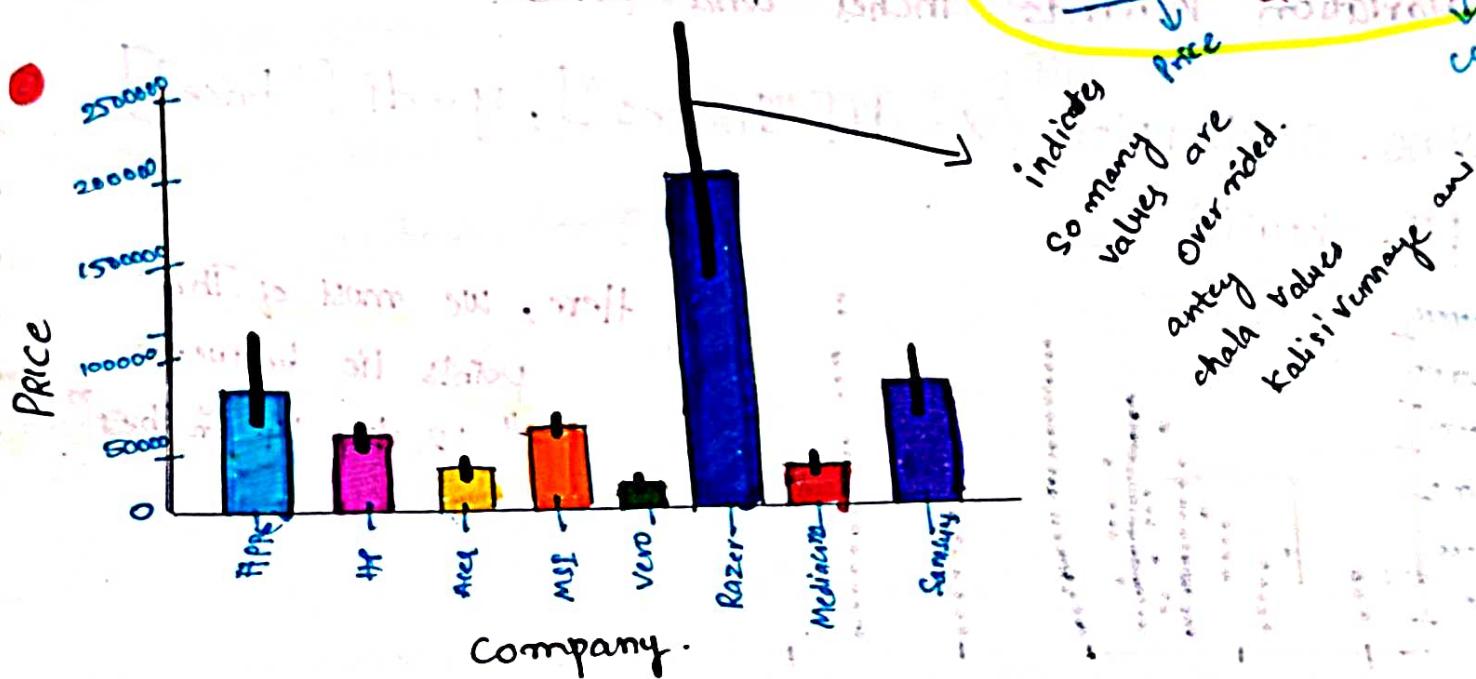


df[["TypeName"]]
 .value_counts()
 Notebook: 727
 Gaming: highest
 2 in 1: 127
 workkit: 29
 Netbook: 25

sns. to
 visual gg
 represent
 chegham

- # average price For each of The Laptop brands
This will say us insights that as per company The price of Laptop vary. \Rightarrow Barplot \Rightarrow discrete Data

- Plt. figure (figsize = (7,3))
sns. barplot ($x = df["Company"]$, $y = df["Price"]$)
plt. xticks (rotation = "vertical")
plt. show()



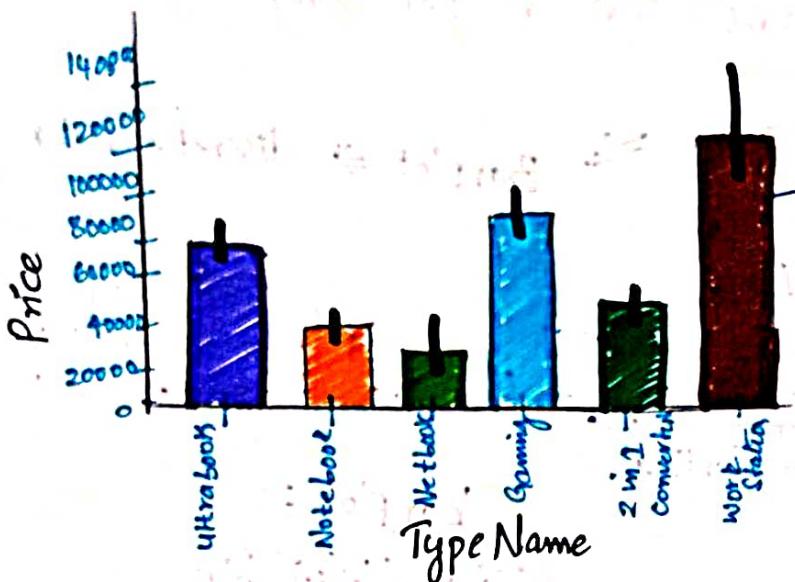
bar plot :-
relationship b/w categorical and data

indicates so many values are overlaid.
antey values chala values kaisi runaye aur

- # laptop type and variation about price

- sns. barplot ($x = df["TypeName"]$, $y = df["Price"]$)
plt. xticks (rotation = "vertical")

Graph.

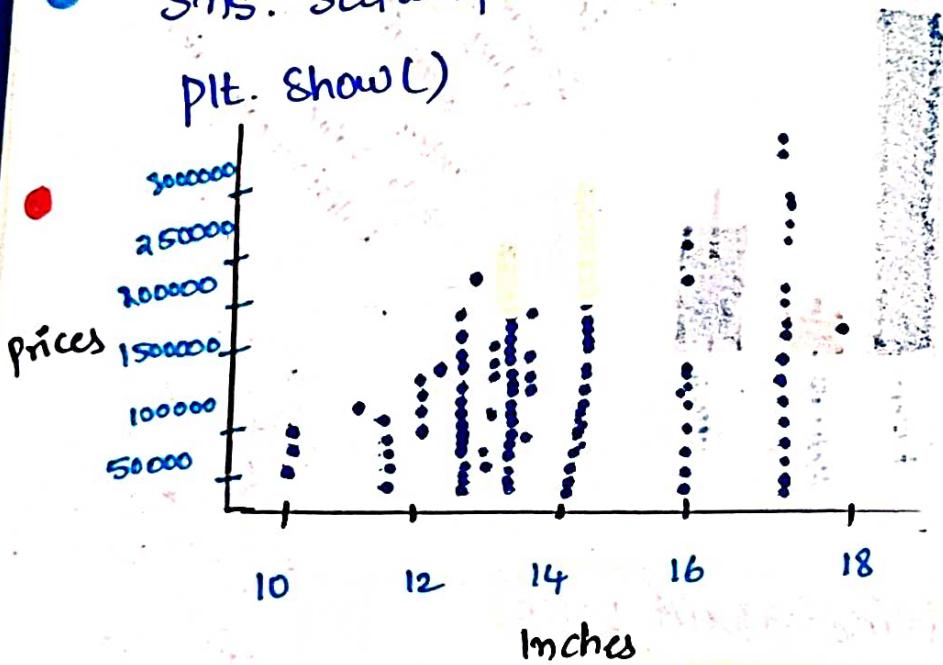


Here, workstation cost or price is more compare to others, because very few buys these laptops. That's why value-counts have 29 only.

- Variation Towards "inches" and "Price".

sns. Scatterplot ($x = df["Inches"]$), $y = df["Price"]$

plt. show()



Here, we most of the points lie between 12 to 14 inches

For screen resolution column we have many types of screen resolutions out there as shown. "Touch Screen" and "Normal" and IPS panel are 3 parts on basis of which we can separate the things.

df[["Screen Resolution"].value_counts()]

Full HD 1920x1080

507

1366x788

281

IPS panel Full HD

230

IPS panel / Touch Screen

53

Here, we see,

repeated words in

Screen Resolution:

Touch screen / Full HD 1920x1080

1

* Full HD

Touch screen / Quad HD x 3200

1

* IPS panel

Touch screen / 4K ultra

1

* Touch screen

Create a New column, Touch Screen if The value is "1" That

Laptop is Touch Screen.

Take no. of arguments.
in one expression.

df[["Touch Screen"]]

New column.

element : 1
if "Touch screen" in element else 0

(lambda element : 1
if "Touch screen" in element else 0)

ardham

Lambda function, For loop la [element = Full HD 1920x1080]

Value_counts lo. eKKada Touch Screen

Vuntadho, akkada "1"

ani create chesthadi.

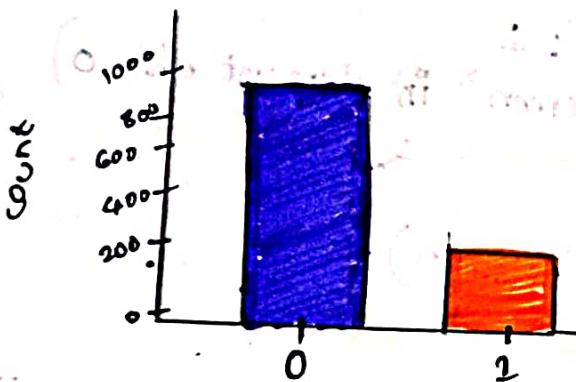
else (lekapothey) "0"

ani df[["Touch Screen"]] = column lo
Save ayethadhi.

EKKADA, TouchScreen, Vunkey akkada \rightarrow 1 lekapothe "0"

Company	Type Name	Inches	Screen Resolution	CPU	Ram	Mem.	Price	Touch Screen
1298	ASUS	15.6	IPS Panel Retina Display	i7	8GB	1TB	Rs. 111111	0
1299	ASUS	15.6	IPS Panel Full HD Touchscreen	i7	8GB	1TB	Rs. 111111	1
1302	ASUS	15.6	IPS Panel Quad/ Touchscreen	i7	8GB	1TB	Rs. 111111	0
			1920x1080					

- df[["Touch Screen"]].value_counts()
 - 0 = 1111
 - 1 = 192
- sns.countplot(df[["Touch Screen"]])

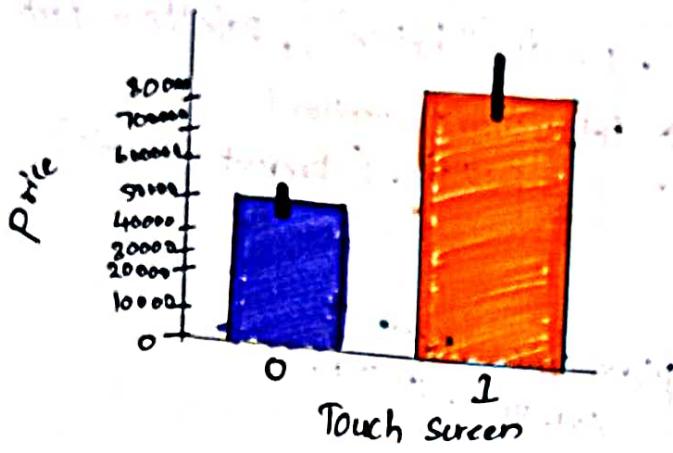


Touch Screen: Comparison with price of laptop

Touch Screen: Comparison with price of laptop

sns.barplot(x=df[["Touch Screen"]], y=df[["price"]])

Graph.



From Here, we can say, The Price for The Touch Screen is more Than normal Desktop ps



`df[["Ips"]] = df[["ScreenResolution"]].apply(lambda`

`element : 1`

`if "Ips" in element else 0)`

`df.Sample(5)`

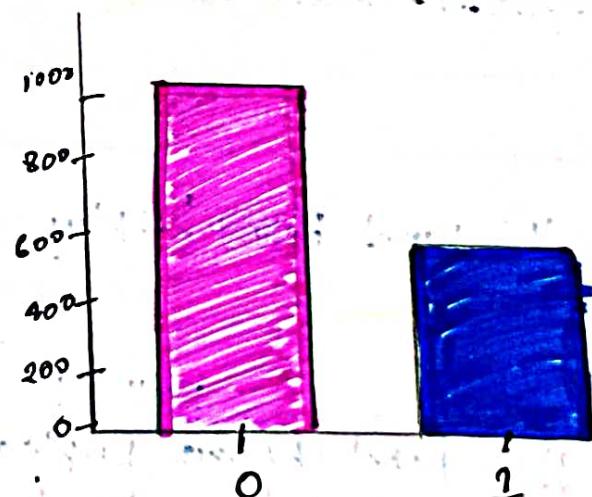
	Company	Type Name	Inches	ScreenResolution	-----	Price	Touch screen	Ips
263				Full HD		0	0	0
396				Ips Panel		0	1	1
751				1366 x 768		0	0	0
221				Ips panel Full HD		0	1	1
112				Full HD		0	0	0

`df[["Ips"]].value_counts()`

0 938

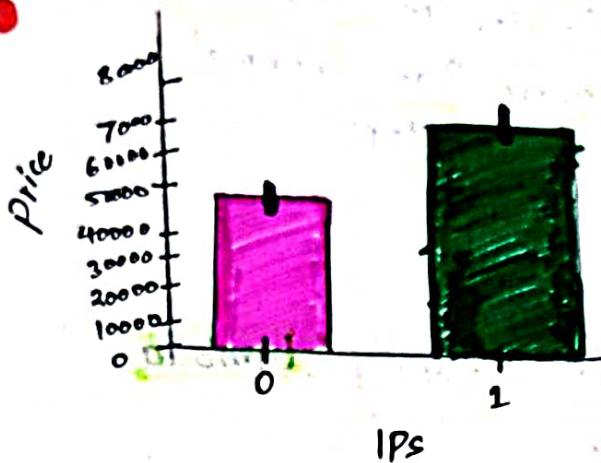
1 365

`sns.Catplot(df[["Ips"]])`



Sns. barplot ($x=df["\text{IPS}"]$, $y=df["\text{Price}"]$, palette="dark")

intrusion prevention system. (prevent threats)



Extracting The X Resolution and Y Resolution.

- # We will Split The Text at the "x" letter and Separate The 2 parts
- # From This we can Observe That One of The col is Y resolution
- # Some Feature Engineering
- splitdf = df[["ScreenResolution"]].str.split("x", n=1, expand=True)

Ex:- 1280×1600
1280 before "x" 1600 after

	0	1
0	IPS panel	Retina Display
1	Balance	Full HD
2	IPS Panel	Retina Display
3	IPS panel	Retina Display

whatever, it may be before "x" = 1440

Take That into consideration. even if "0" also

balance Space Consider

Split df = df["Screen Resolution"].str.split("x", n=1, expand=True)

fit

df["X_res"] = split_df[0]

df["Y_res"] = split_df[1]

df.head()

Extract column.

Company	Type Name	inches	Screen Resolution	Price	Touch screen	IPS	X-res	Y-res
0			IPS Panel Retina Display 2560 x 1600				IPS Panel 2560	1600
1			1440 x 900				1440	900
2			FHD 1920 x 1080				FHD 1920	1080
3			IPS Panel Retina 2880 x 1800				IPS 2880	1800
4			IPS Panel Retina 2560 x 1600				IPS 2560	1600

So, basically from that whole text of the "X-res" col, we need to extract the digits from it, but the problem is the numbers are scattered in some cases, that is the reason why I am trying to use regex, if we use this we will exactly get the numbers which we are looking for. So, firstly replace all the numbers from that string as "Id+?". Then find all numbers.

|d means That integer

|.? all numbers which come after a number

|d+ The string must end with number

• EKKada manam. only numbers consider chesi, Text ignore cheesam

• $df['x_res'] = df['x_res'].str.replace(',', '')$

• str.findall(`(\d+(\.\d+)?)`).apply^{map}

Sorting
"X" to $x[0]$
any Regular Expression [RegEx]
first dhammeda apply chetham
ani ardhama.
and ends with integer
There is (.) dot not in The text data.

EX: IPS Panel Retina Display

First, we look for number $(\d+)(\.\d+)?(\d+)$

Started with number 2880 ended with integer number

df.head()

Company	Type	Inches	Screen Resolution	Price	Touch Screen	IPS	X-res	Y-res
ASUS	Laptop	15.6	IPS Panel Retina display 2560x1600	1440	Yes	Yes	2560	1600
ASUS	Laptop	15.6	IPS Panel Retina display 2560x1600	1440	Yes	Yes	1440	900
ASUS	Laptop	15.6	IPS Panel Retina display 2560x1600	1920	Yes	Yes	1920	1080
ASUS	Laptop	15.6	IPS Panel Retina display 2560x1600	2880	Yes	Yes	2880	1800
ASUS	Laptop	15.6	IPS Panel Retina display 2560x1600	2560	Yes	Yes	2560	1600

Ela mancham.

Ela mancham.

Convert $[x_res][y_res]$ column into "int" Data

df["x_res"] = df["x_res"].astype("int")

df["y_res"] = df["y_res"].astype("int")

df["y_res"] = df["y_res"].astype("int")

df.info()

Data columns

Dtype

object

On Company

13 x-res

int32

14 y-res

int32

old : object

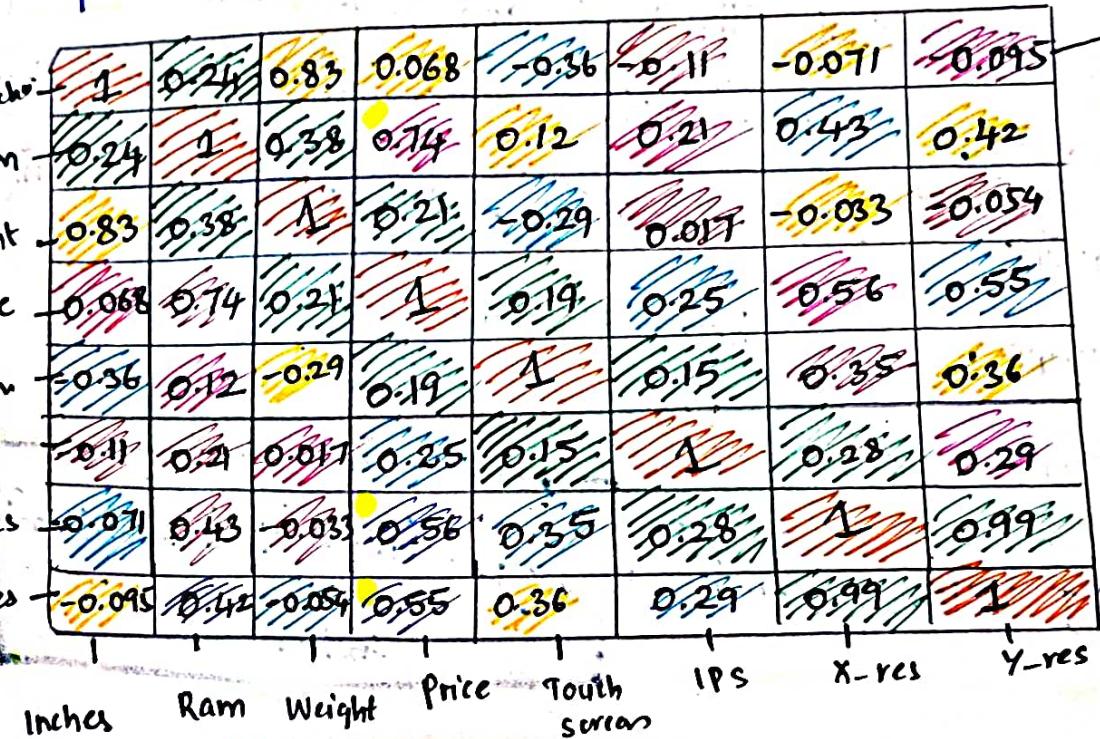
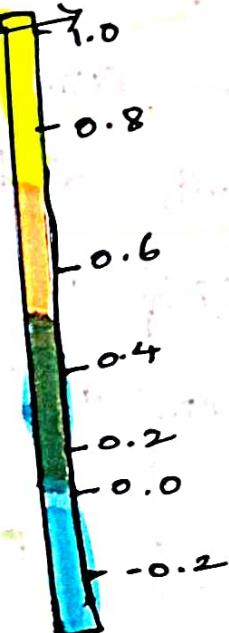
old : object

Correlation

sns. heatmap (df.corr(), annot=True, "plasma")

plt. figure (fig_size=(15,7))

Display number



df.corr() ["Price"]

Inches 0.068197

Ram 0.743007

Weight 0.20370

price 1.000000

Touch screen 0.191226

IPS 0.252208

x-res 0.556529
y-res 0.552809

From The correlation plot we observed that the X_{res} and Y_{res} is increasing, The price of the laptop is also increasing, so X_{res} & Y_{res} are positively correlated and They are giving much information so that is the reason. why I had Spitted resolution column into X_{res} & Y_{res} columns respectively.

So, to make things good, we can create a new column named PPI [Pixel Per Inch], Now as we saw from correlation plot that X_{res} and Y_{res} are having much collinearity, so why not combine them with inches. which is having "less collinearity", so, we will combine them as follows, so here formula

To Calculate PPI [Pixel Per Inch]

$$PPI [\text{Pixel per inch}] = \frac{\sqrt{X_{res} \text{ Solution}^2 + Y_{resolution}^2}}{\text{Inches}}$$

To combine X_{res} , Y_{res} and Inches. They must have same datatypes

$$* \text{Inches} = \text{Float64} \rightarrow$$

$$* X_{res} = \text{Int}$$

$$* Y_{res} = \text{int}$$

$$df['PPI'] = \left(\frac{\sqrt{(df['x-res'])^2 + df['y-res'])^2}}{df['Inches']} \right)^{0.5}$$

PPI = $\frac{\sqrt{(x\text{-res})^2 + (y\text{-res})^2}}{\text{Inches}}$

Divide by

`df.head()`

Company	Type	Plan	Inches	Screen Resolution	X-res	X-res	PPI
0			13.3		2560	1600	226.983
1			13.3		440	900	127.677
2			15.6		1920	1080	141.21199
3			15.4		2880	1800	220.5346
4			13.3		2560	1600	226.98300

EX:-
$$\frac{(2560 + 1600)^{0.5}}{13.3} = 226.983$$

`df.corr()["Price"]`

0.0681	0.556529	0.552809	0.473487
Inches	X-res	Y-res	PPI

We combined these three columns with which gives lots of Good correlation.

So, we drop These 3 columns along with "screen resolution" column because we extracted total data from this column.

and also we can see that there is high correlation between X-res and Y-res so we can drop one of them up above.

`df.drop(columns = ["ScreenResolution", "Inches", "x-res", "y-res"], inplace = True)`

`df.head()`

Company	Type Name	Cpu	Ram	Memory	Gpu	Opsys	Weight	Price	Touch Screen	IPS	PP2

↳ Remaining Columns-

We will Work on Cpu Column

`df["Cpu"].value_counts()`

Intel Core i5 7200U 2.5GHz 190
Intel Core i7 7700HQ 2.8GHz 146
Intel core i7 7500U 2.7GHz 134
....

AMD A6-Series

Most Common Processors are made by Intel right, so we will be clustering their processor into different categories. i5, i7, Other → Processors of Intel which do not have i3, i5, i7 attached to it. Completely different so, that's the reason i will cluster them into Other and Other category is "AMD" which is different category in whole. So, if we observe we need to extract the first 3 words of CPU column as first words of every row under CPU column of the CPU, we will use them.

- Text = "Intel Core i5 7200U 2.5GHz"
 Text.split() (split chesinam)
 ["Intel", "Core", "i5", "7200U", "2.5GHz"]
- Text.split()[3] (It gives only upto 2)
 ["Intel", "Core", "i5"] (dhani only first three matramay select) (Breaking code into parts)
- " ".join(Text.split()[3])
 "Intel core i5" Tarwata join chesam.
- " ".join(Text.split()[3])
 "Intel Core i5" → dhani ardham.: space echam.
- df["cpu-name"] = df["Cpu"].apply(lambda text: " ".join(text.split()[3]))

Company	Type name	Cpu	Cpu-name
Lenovo		intel Core i5 2.3GHz Intel core i5 1.8GHz ... intel core i5 3.1 GHz	intel Core i5
HP			intel Core i5
Dell			intel Core i5
Asus			intel Core i5

Problem Here : Some rows have in Cpu-name

Have [ss8.intel Atom x5-Z8350
204.intel Xeon E3-1505M]

So, Manaki only i5, i3, i7
Vunnaye. Kakapotheey. EKKada
ayethey Vati simple ga "Intel" rasi, migithaye AMD Processor
matramay CKKuva "ur intel" lo kuda rakalu vun
anigadham.

Function:

def processor_type(text):

if text == "Intel Core i7" or text == "Intel Core i5"
or text == "Intel Core i3":

return text

else:

if text.split()[0] == "Intel":

return "Other Intel processor"

else:

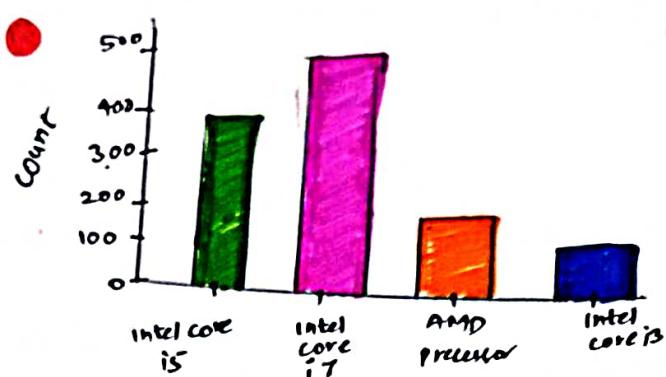
return "AMD processor"

df[["cpu-name"]] = df[["cpu-name"]].apply(lambda text:
processor_type(text))

df. Sample(5)

Company	Type Name	CPU	Processor	CPU-name
		Intel Celeron Dual core		AMD processor
		AMD A9 Series		AMD Processor
		Intel Core i7 7700 HQ		intel core i7
		Intel Celeron Dual core		AMD Processor
		Intel Core i7		intel core i7

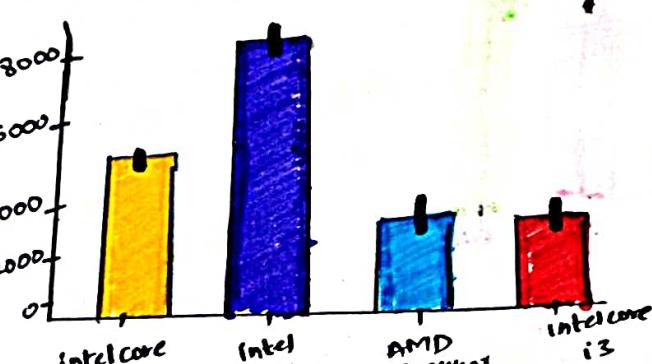
Sns. countplot (df[["cpu-name"]])



remove 'Other processor'

issue [other processor]
any column was selected ▲

Variation between "Price" & "cpu-name"



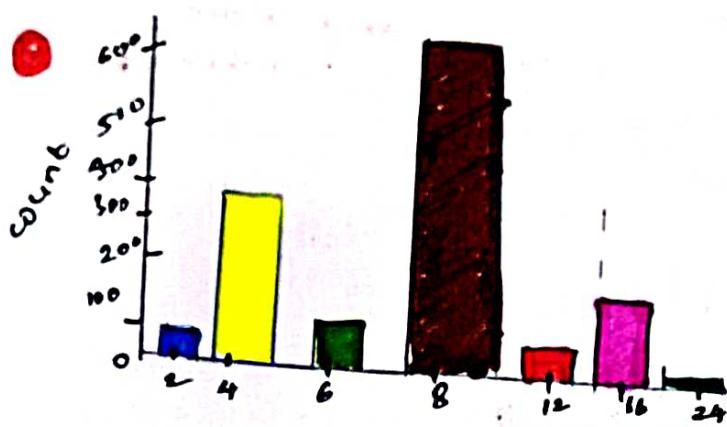
df.drop(columns=[\"cpu\"], inplace=True)

df.head()

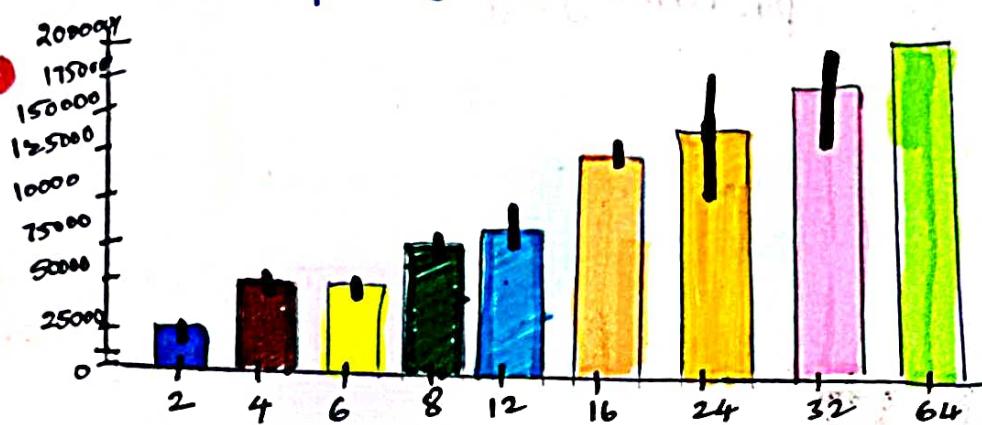
Company	Type Name	Ram	Memory	Gpu	Opsys	Weight	Price	Touch Screen	IPS	PPI	CPU
ASUS	ASUS ROG Strix G15	16GB	1TB	NVIDIA GeForce RTX 3060	Windows 11	2.2 kg	₹ 65,990	Yes	15.6	250	i5
ASUS	ASUS ROG Strix G15	16GB	1TB	NVIDIA GeForce RTX 3060	Windows 11	2.2 kg	₹ 65,990	Yes	15.6	250	i5
ASUS	ASUS ROG Strix G15	16GB	1TB	NVIDIA GeForce RTX 3060	Windows 11	2.2 kg	₹ 65,990	Yes	15.6	250	i5
ASUS	ASUS ROG Strix G15	16GB	1TB	NVIDIA GeForce RTX 3060	Windows 11	2.2 kg	₹ 65,990	Yes	15.6	250	i5

RAM Column

- sns.countplot(df["Ram"])



- sns.barplot(x=df["Ram"], y=df["Price"])



Memory Column

We will Separate the type of memory and the value of it, just similar to one done in previous part. This part involves various steps, we have it in different dimension as a Complete 128GB SSD + 1TB HDD, so in order to for it come in same dimension, we need to do some modification.

We have different categories and also different variations

df["Memory"].value_counts()

256 GB SSD

412

1 TB HDD

223

128 GB SSD + 1 TB HDD

94

1.0 TB Hybrid

9

128 GB Flash Storage

4

HDD, SSD, Flash, Hybrid

We have 4 common variants

for ex: 1.0 TB will be

This Expression

remove

decimal space

1TB

df["Memory"] = df["Memory"].astype(str).replace

Ex:- 128GB SSD + 1TB HDD

("1.0", "", regex=True)

what ever coming
after decimal point
Delete it.

Ex:- 1.0 TB = 1TB

convert it in
String.

128 SSD + 1 TB HDD

128 SSD + 1000 HDD

replace GB with "

apply to
Every tuple

df["Memory"] = df["Memory"].str.replace("GB", "")

replace TB with "000"

df["Memory"] = df["Memory"].str.replace("TB", "000")

split the word across "+" character

newdf = df["Memory"].str.split("+", n=1, expand=True)

128 SSD + 1000 HDD

["128 SSD", "1000 HDD"]

Removed

new df

0	1
128 SSD	None
128 Flash storage	None
256 SSD	None
512 SSD	None
⋮	
1000 HDD	None
500 HDD	None

newdf. Sample(5)

0	1
"	None
"	1000 HDD
"	None
"	None
1000 SSD	2000 HDD

Endnote
nude values
management

1000 SSD + 2000 HDD

we will strip up all the "white spaces", basically Eliminating white space.

df["first"] = newdf[0] # here we considered only [0], but ignore [1]

df["first"] = df["first"].str.strip()

df.head()

string ga \Rightarrow strip()
marchi apply changes
white space notation
portage.

Company	Type Name	Inches	Memory	first
Enduro			1000 HDD	1000 HDD
Weka			512 SSD	512 SSD
Winagle			1000 HDD	1000 HDD
Hybrid			256 HDD	256 HDD
Flash Storage			512 HDD	512 HDD

↓ white
with space

without space

def apply changes (values):

df["Layer 1" + value] = df["first"].apply(lambda x: 1 if value in x else 0)

List to apply = ["HDD", "SSD", "Hybrid", "Flash Storage"]

for value in List to apply :

apply changes (values)

Detail explanation

`def apply changes (value)`

important
 $\text{df}[\text{"Layer1"} + \text{Value}] = \text{df}[\text{"first"}].apply[\lambda x: 1 \text{ if Value in } x \text{ else 0}]$

$\text{Layer1} + 1000\text{HDD}$
 Vunna chaggarai "value" available

$\text{List to apply} = [\text{"HDD"}, \text{"SSD"}, \text{"Hybrid"}, \text{"Flash"}]$ storage

$\text{for Value in List.to apply:}$
 $\text{apply changes (values)}$

def Function

Lambda function dheni meda apply cheyali

`df.head()`

Company	Type Name	Ram	Memory	first	Layer1 HDD	Layer1 SSD	Layer1 Hybrid	Layer1 Flash Storage
Company 1	SSD	128	SSD	128 SSD	0	1	0	0
Company 2	SSD	128	SSD	128 Flash Storage	0	0	0	1
Company 3	SSD	1000	HDD	1000 HDD	1	0	0	0
Company 4	SSD	128	Hybrid	128 Hybrid	0	0	1	0

remove all characters just keep the numbers

`df["first"] = df["first"].str.replace("128 SSD", "128")`

`df["first"].value_counts()`

$256 \rightarrow 508$
 $1000 \rightarrow 250$

$128 \rightarrow 177$	$508 \rightarrow 1$
$512 \rightarrow 140$	$240 \rightarrow 1$
$8 \rightarrow 1$	

remove all text

Ex:- 128 SSD

A:- 128

we consider **Second Part** of newdf[0][1].

1200 SSD
records

1200 SSD + 120 HDD
records

df["Second"] = newdf[1]

df.head()

Company	Type Name	inches	Processor	First	Layers 1 SSD	Layers 1 HDD	Layers 2 Hybrid	Flash	Second
Apple	MacBook Pro	13	i5	128	128	0	0	0	None
Apple	MacBook Pro	15	i7	128	128	0	0	0	1000 HDD
ASUS	Transformer	13	i5	256	0	0	0	0	None
ASUS	Transformer	13	i7	512	0	0	0	0	None
ASUS	Transformer	13	i7	256	0	0	0	0	2000 HDD

We have apply same as [0].

def apply changes 1 [value]:

df["Layer 2" + Value] = df["Second"].apply(lambda x: 1 if Value in x else 0)

List to apply 1 ["HDD", "SSD", "Hybrid", "Flashstorage"]

df["second"] = df["Second"].fillna("0")

for value in list to apply 1:

apply changes 1 [value]

df["Second"] = df["Second"].str.replace("/", "D,")

df["Second"].value_counts()

0 → 1095

1000 → 187

256 → 3

2000 → 15

500 → 2

512 → 1

change datatypes of "first" and "second" to int

• $df["first"] = df["first"].astype("int")$
 $df["Second"] = df["Second"].astype("int")$

$df.tail()$

$df.info()$

note: first and second are object type, the others are int

first
:
second

int : old: first \Rightarrow object
int : old: second \Rightarrow object

multiplying the elements and storing the result in

subsequent columns

• $df["HDD"] = (df["first"] * df["layer1 HDD"] + df["second"] * df["layer2 HDD"])$

Ex:- $(128 \times 0) + [0 \times 0] = 0$

• $df["HDD"] = 0$

$$[1000 \times 1] + [0 \times 0] = 1000$$

• $df["HDD"] = 1000$

$$df["SSD"] = (df["first"] * df["layer1SSD"]) + \\ df["Second"] * df["layer2SSD"]$$

$$df["Hybrid"] = (df["first"] * df["layer1Hybrid"]) + \\ df["Second"] * df["layer2SSD"]$$

$$df["Flash-storage"] = df["first"] * df["layer1Flash-storage"] \\ + df["Second"] * df["layer2Flash-storage"]$$

dropping of unnecessary columns

```
df.drop(columns = ["first", "Second", "layer1HDD",
                    "layer1SSD", "layer1Hybrid", "layer1Flash-
                    storage", "layer2HDD", "layer2SSD",
                    "layer2Flashstorage",
                    "layer2Hybrid"], inplace = True)
```

f. Sample (5)

$$\text{Cost} = [0 \times 0] + [1 \times 0001]$$

$$\text{Cost} = [0 \times 0011] \text{ fb}$$

comp_id	Type	Non_Ram	Ram	Memory	...	HDD	SSD	Hybrid	Flash_Storage
27	Dell	Notebook	8	256SSD		0	256	0	0
680	MSI	Gaming	8	256SSD		0	256	0	0
431	Lenovo	Notebook	8	256SSD		0	256	0	0
907	Acer	Notebook	4	16 Flash storage		0	0	0	0
461	Dell	Notebook	8	2000HDD		2000	0	0	0

df.drop(columns = ["Memory"], inplace = True)

df.head()

column will be deleted

(malla Evadu Rayali ☺)

df.corr().Price

Ram 0.743

Weight 0.210

Price 1.000

Touchscreen 0.191

IPS 0.252

PPI 0.473

HDD -0.096

SSD 0.670 ✓

Hybrid 0.007

Flash_Storage NaN

* not contributing to price

column. (or) very less contributing.

* We can drop these two columns.

Hybrid & Flash_Storage.

* where SSD & HDD having good correlation

* But, HDD has (-ve) negative relation with price, there is more probability that

* The price of Laptop is increasing there is more use of SSD instead of HDD and vice versa

df.drop(columns = ["Hybrid", "Flash_storage"], inplace = True)

Company	Type	Ram	Gpu	opsys	Weight	Price	Touch	IPS	PPI	CPU	HDD	SSD

GPU column.

df.Gpu.value_counts()

intel HD Graphics 620 281

intel HD Graphics 520 185

⋮ ⋮

ARM Mali T860 MP4 1

Here as we are having less data regarding the laptops, it's better that we focus on GPU brands instead focusing on the values which are present there beside them, we will focus on the brands.

This is what we will be doing, Extracting the brands.

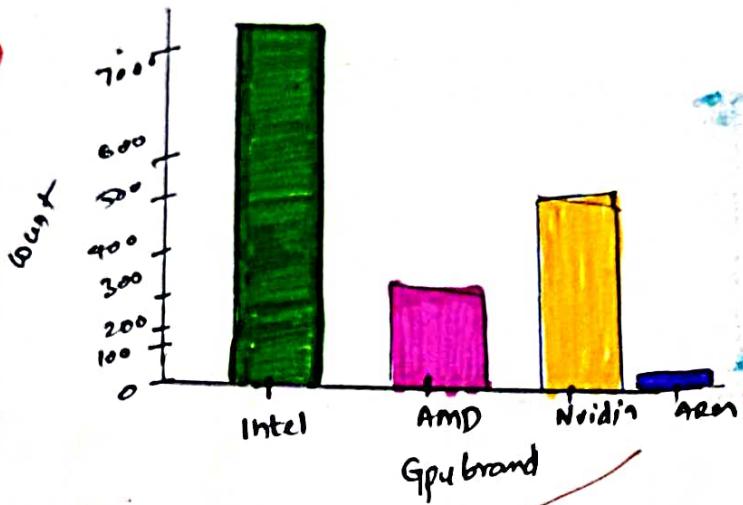
a = df["Gpu"].iloc[0]

Index(0)

Print(a.split()[0]).

Intel

df[["Gpu brand"]] = df[["Gpu"]].apply(lambda x: x.split()[0])
sns.countplot(df[["Gpu-brand"]], palette="plasma")

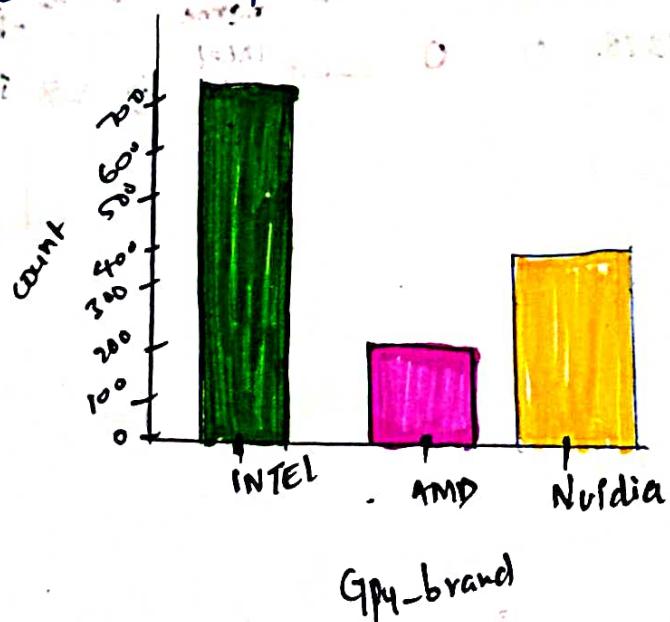


Remove "ARM" tuple

df = df[df[["Gpu-brand"]] != "ARM"]

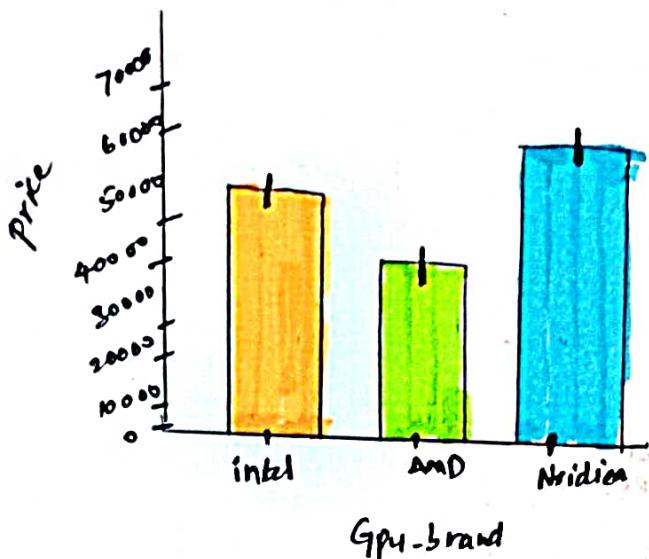
df = df[df[["Gpu-brand"]]]

sns.countplot(df[["Gpu-brand"]])



price - Gpu analysis, i used np.median in order to
check if there is any impact of outlier or not

- sns.barplot (df[["Gpu-brand"]], df[["Price"]], estimator = np.mean).



- $df = df.drop([["Gpu"]], columns=)$

df.head

company	typeName	ram	opsys	weight	price	toucScreen	ips	PPI	CPU_name	HDD	SSD	Gpu -brand
Apple	ultra book	8	macOS	1.37	71378.	0	0	226.9	Intel core i5	0	128	intel

OpSys column

- $df[["Opsys"]].value_counts()$

Windows 10 1072

No OS 66

Linux 62

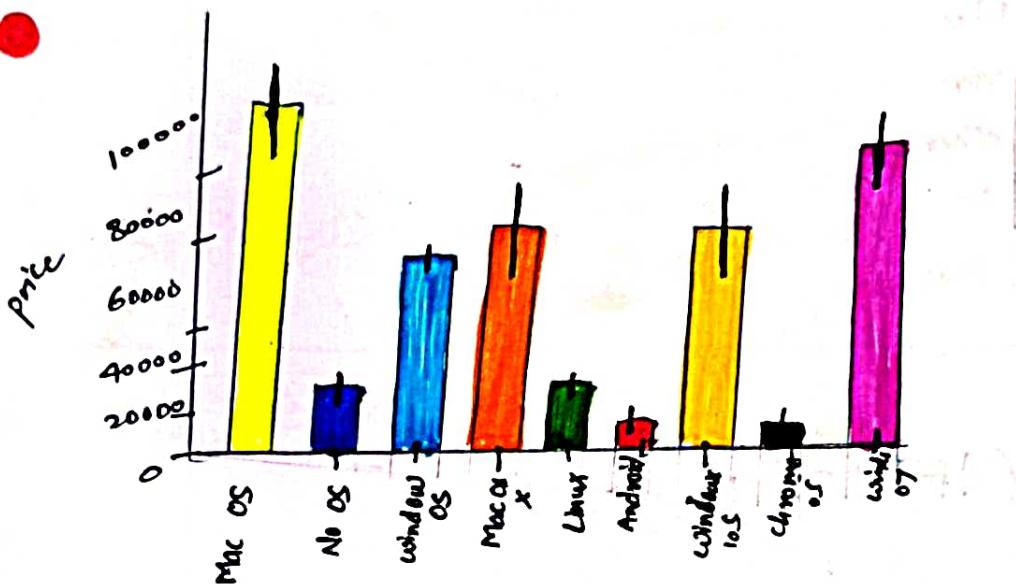
Android 8

Variation factor b/w "price" and "opsys" column.

• sns.barplot(df["opsys"], df["Price"])

plt.xticks(rotation="vertical")

plt.show()



We apply same strategy as Memory column.

def SetCategory(text):

if text == "windows 10" or text == "windows":
 or text == "windows 10 s":
 return "windows"

elif text == "Mac OS X" or text == "Mac":
 return "Mac"

else:

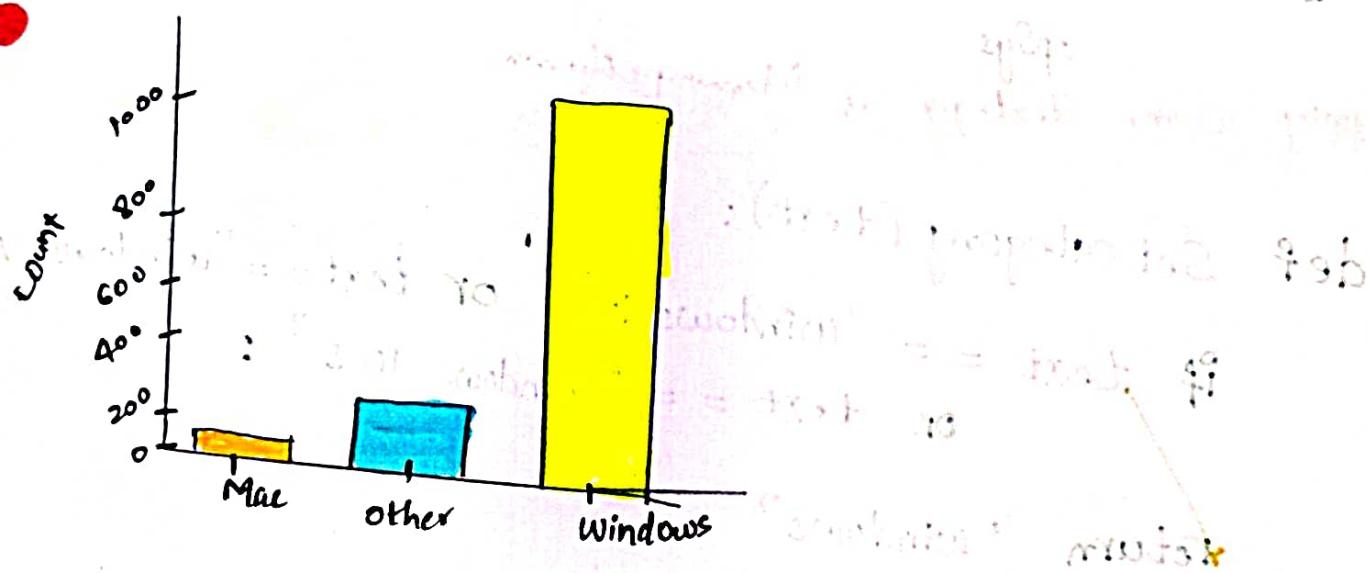
 return "other"

`df[["opSys"]] = df[["OpSys"]].apply [lambda x: set category(x)]`

`df.head(5)`

Company	RfpName	Ram	OpSys	Processor	GPU	Price	GpuLam
Apple			Mac				
Apple			Mac				
HP			Other				
Apple			Mac				
Apple			Mac				

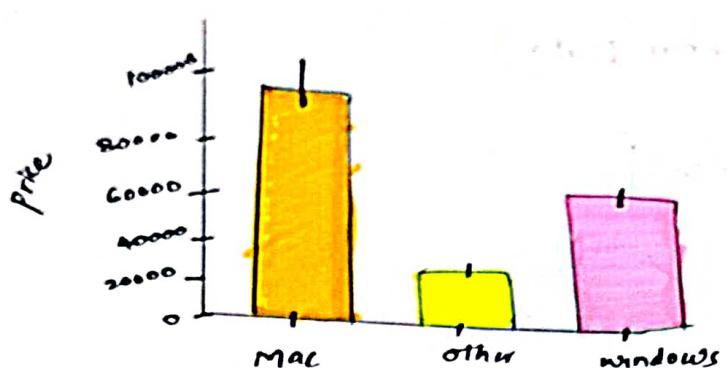
`sns.countplot(df[["OpSys"]])`



Variation between "OpSys" & "Price"

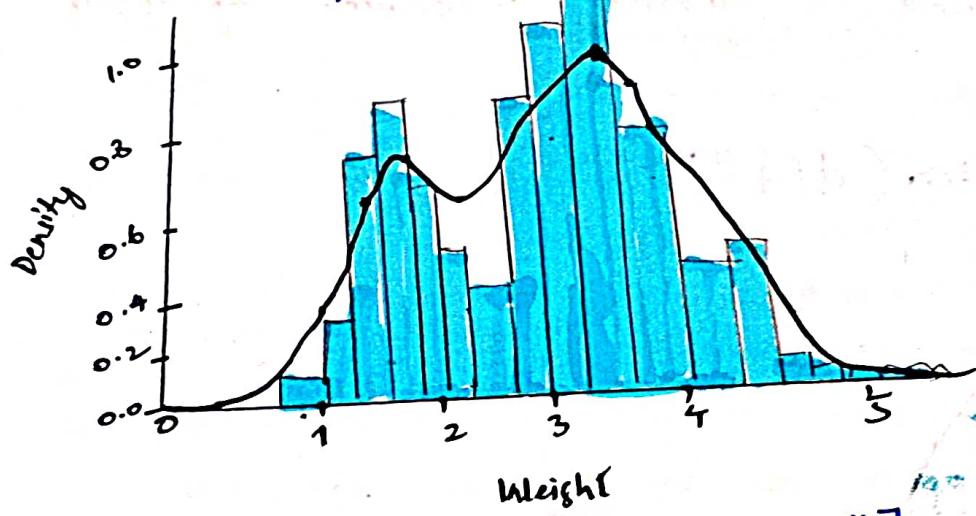
`sns.barplot(x=df[["OpSys"]], y=df[["Price"]])`

`plt.xticks(rotation = "vertical")`

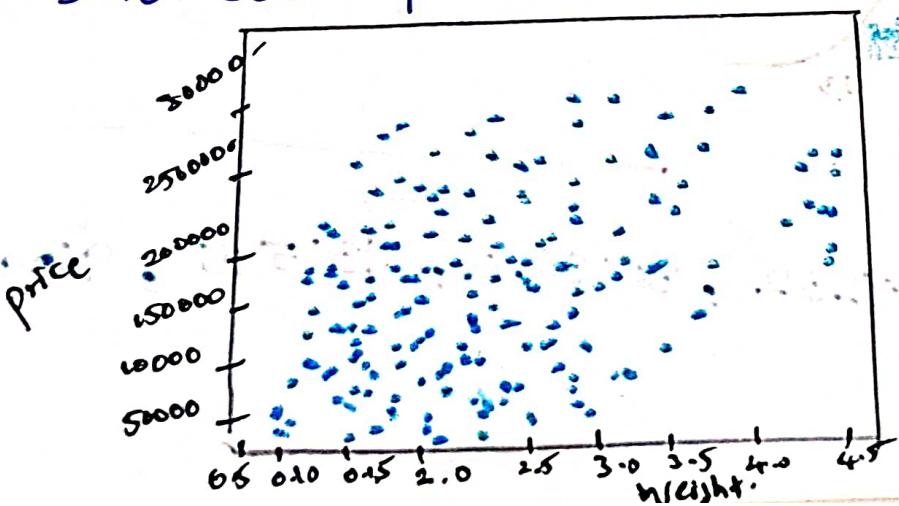


Weight Analysis

`sns.distplot(df[["weight"]])`

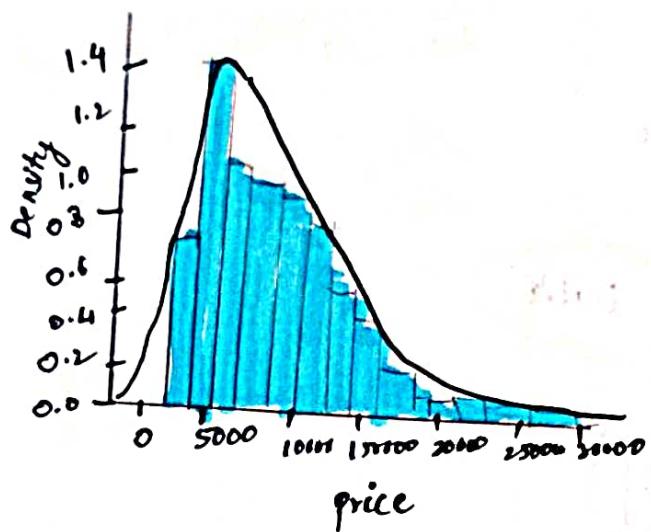


`sns.scatterplot(df[["weight"]], df[["price"]])`



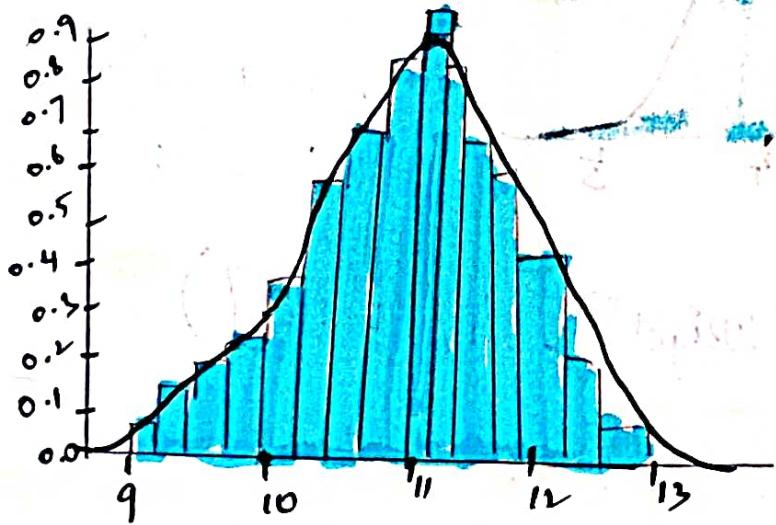
Price Analysis → Continuous Data

- sns.distplot(df["price"])



we convert right skewed To Normal distribution by apply "log" to it.

- sns.distplot(np.log(df["Price"]))



- list(df["Price"])[::5]

[71378.6382, 47895.5232, 30636.0, 131595, 96095.8]

correlation with price.

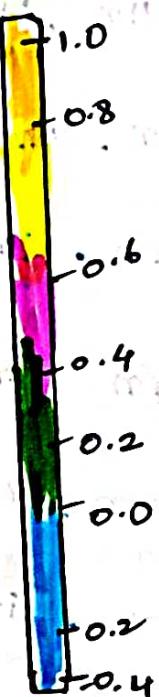
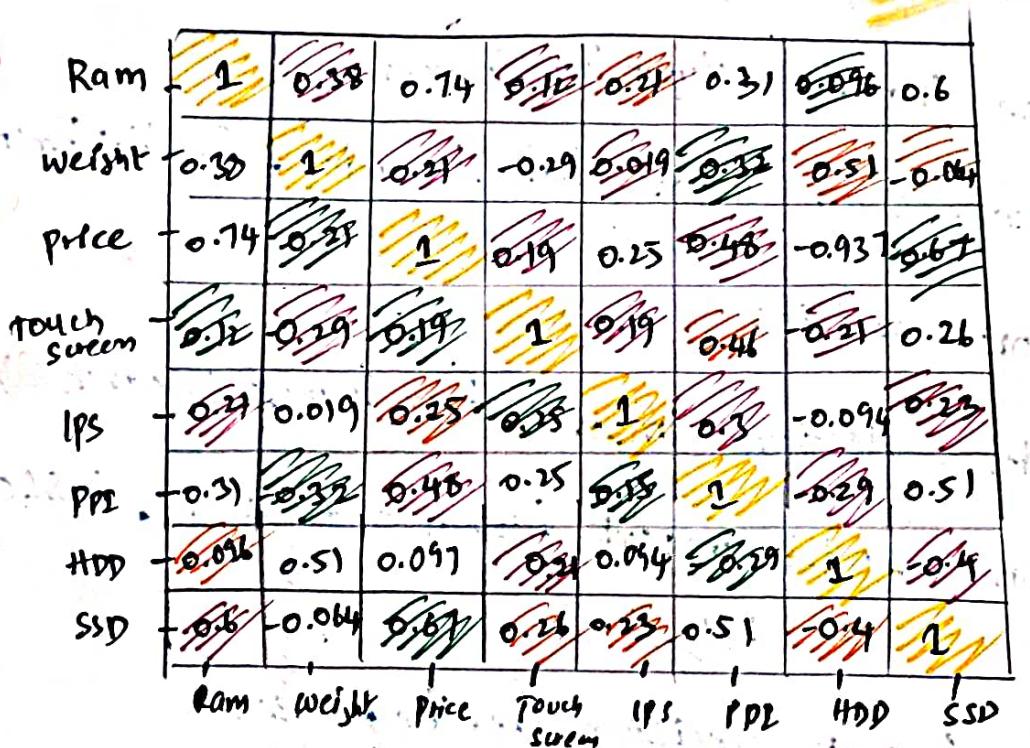
df.corr()["price"]

Ram	0.14290
Weight	0.2098
Price	1.0000
Touchscreen	0.192917

IPS	-0.253320
PPI	-0.475
HDD	-0.0968
SSD	-0.670660

plt.figure(figsize=(10,5))

sns. heatmap (df.corr()), annot=True, cmap="plasma")



continuation of converting "log" of price

np.log(71378.6832)

11.175754549

np.exp(11.1757549)

71378.6832

MODEL BUILDING

```
test = np.log(df["Price"])
train = df.drop(["Price"], axis=1)

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
from sklearn import metrics
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LinearRegression, Lasso, Ridge
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from xgboost import XGBRegressor
from sklearn.svm import SVR
from sklearn.neighbors import KNeighborsRegressor
```

```
from sklearn import tree
```

Train Test Split

- X-train, X-test, y-train, y-test = train-test-split
(train, test, test_size = 0.15,
random_state = 2)

- X-train.shape, X-test.shape
((1106, 12), (196, 12))

There's a class which we imported named as ColumnTrans.
We use this widely while building our models using
former
Pipelines
so for this we have to get the index numbers
of the columns which having categorical variables.
mapper = {
 i: Value for i, Value in enumerate
 (X-train.columns)}

mapper

0	"company"	(Categorical)	6	"Ips"
1	"Type Name"	(cat)	7	"PPJ"
2	"Ram"		8	"cpu-name"
3	"opSys"	(cat)	9	"HDD"
4	"Weight"		10	"SSD"
5	"TouchScreen"		11	"Gpu-brand"

Linear Regression

- # Here, we apply OneHotEncoding for columns :- 0, 1, 3, 11
- # The Remainder we keep as "pass through", i.e. no other col must get effected.
- # except the ones undergoing transformation.
- # → Converting Categorical variable To Numeric Variable
- Step 1 = Column Transformer (transformers = [("col_tnf",

OneHotEncoder (sparse = False, drop = "first"), [0, 1, 3, 8, 11]), remainder = "pass through")
handle_unknown = "ignore")
we have lots of "0" in
The matrixes
Ex: [1, 0, 0, 0, 0, 0, 0, 0]
An identity matrix of dimension.
[1000x1000] diagonals are 1.
rest all are zero.

model
Step 2 = LinearRegression()
label (ColumnTransformer)
Pipe = Pipeline([("step1", step1), ("step2", step2)])

Pipe.fit(x-train, y-train)
fitting the data

A predict

$$y_pred = \text{Pipe}.predict(x_test)$$

Point("R2 score", metrics.r2_score(y_test, y_pred))

.Point("MAE", metrics.mean_absolute_error(ytest, yPred))

$$\begin{array}{ll} \text{R2 score} & 0.8073 \\ \text{MAE} & 0.21017 \end{array}$$

Giving error %
0.21%

$$\text{np.exp}(0.21)$$

$$1.2336780$$

$$71378.6832 - \text{np.exp}(0.21017)$$

$$71377.449301$$

it is giving 1.233 difference.

Ridge Regression

- Step 1 = ColumnTransformer (transformers = [`[{"cat-trf":`
`OneHotEncoder (sparse = False, drop = "first",`
`handle_Unknown = "Ignore")], [0, 1, 3, 8, 11]]`
`remainder = "passthrough"`])
 - Step 2 = Ridge (alpha = 10)
 - Pipe = pipeline([("step1", step1), ("step2", step2)])
 - Pipe.fit(x_train, y_train) # fitting data
 - y-pred = pipe.predict(x-test) → on test data
 - print("R2 Score", metrics.r2_score(y-test, y-pred))
 - print("MAE", metrics.mean_squared_error(y-test, y-pred))
 - R2-score 0.8120403
 - MAE 0.21034714
- Predicting on "price" column*

Lasso Regression

Step 1 = ColumnTransformer (transformers = [("col_tnf", OneHotEncoder (sparse = False, drop = "first", handle_unknown = "ignore"), [0, 1, 3, 8, 11]), remainder = "passthrough")

Step 2 = Lasso (alpha = 0.01)

Pipe = pipeline[["step1", step1], ("step2", step2)]

y_Pred = pipe.Predict(x-test)

print("R2-score", Metrics.r2_score(y-test, y-Pred))

print("MAE", Metrics.mean_squared_error(y-test, y-Pred))

R2 score 0.80733

MAE 0.2101708

Decision Tree Regressor

Step 1 = ColumnTransformer (transformer = ([{"col_tnf":
OnehotEncoder (sparse = False, drop = "first",
handle_error = "ignore"), [0, 1, 3, 8, 11]],
remainder = "Pass through")])

Step 2 = DecisionTreeRegressor (max_depth = 8)

- Pipe = Pipeline ("step1", step1), ("step2", step2)
- Pipe . fit (x_train, y_train)
- y_Pred = Pipe . predict (x_test)
- Print ("R2_score", metrics . R2_score (y_Pred, y_test))
- Print ("MAE", metrics . mean_absolute_error (y_Pred, y_test)).

R2_Score : 0.848598

mean_absolute_Error : 0.1781070

gulf
26/03/23