

# **Data Science II with python (Class notes)**

**STAT 303-2**

Arvind Krishna

2023-01-03

# Table of contents

<b>Preface</b>	<b>5</b>
<b>I Linear regression</b>	<b>6</b>
<b>1 Simple Linear Regression</b>	<b>7</b>
1.1 Simple Linear Regression . . . . .	7
1.1.1 Training with <code>statsmodels</code> . . . . .	8
1.1.2 Training with <code>sklearn</code> . . . . .	15
1.1.3 Training with <code>statsmodels.api</code> . . . . .	17
<b>2 Multiple Linear Regression</b>	<b>19</b>
2.1 Multiple Linear Regression . . . . .	19
<b>II Assignments</b>	<b>24</b>
<b>3 Assignment 1 (Section 20)</b>	<b>25</b>
Instructions . . . . .	25
3.1 1) Case Studies: Regression vs Classification and Prediction vs Inference (16 points) . . . . .	26
3.1.1 1a) . . . . .	26
3.1.2 1b) . . . . .	26
3.1.3 1c) . . . . .	26
3.1.4 1d) . . . . .	26
3.1.5 1e) . . . . .	27
3.2 2) Examples for Different Regression Metrics: RMSE vs MAE (8 points) . . . .	27
3.2.1 2a) . . . . .	27
3.2.2 2b) . . . . .	27
3.3 3) Modeling the Petrol Consumption in U.S. States (61 points) . . . . .	27
3.3.1 3a) . . . . .	28
3.3.2 3b) . . . . .	28
3.3.3 3c) . . . . .	28
3.3.4 3d) . . . . .	28
3.3.5 3e) . . . . .	28
3.3.6 3f) . . . . .	28

3.3.7	3g)	28
3.3.8	3h)	29
3.3.9	3i)	29
3.3.10	3j)	29
3.3.11	3k)	29
3.3.12	3l)	29
3.4	4) Reproducing the Results with Scikit-Learn (15 points)	30
3.4.1	4a)	30
3.4.2	4b)	30
3.5	5) Bonus Question (15 points)	30
3.5.1	5a)	30
3.5.2	5b)	30
3.5.3	5c)	31
<b>4</b>	<b>Assignment 1 (Sections 21 &amp; 22)</b>	<b>32</b>
	Instructions	32
4.1	1) Case Studies: Regression vs Classification and Prediction vs Inference (16 points)	33
4.1.1	1a)	33
4.1.2	1b)	33
4.1.3	1c)	33
4.1.4	1d)	33
4.2	2) Examples for Different Regression Metrics: RMSE vs MAE (8 points)	34
4.2.1	2a)	34
4.2.2	2b)	34
4.3	3) Simple Linear Regression: Formulation (3 points)	34
4.4	4) Modeling the Petrol Consumption in U.S. States (58 points)	34
4.4.1	4a)	35
4.4.2	4b)	35
4.4.3	4c)	35
4.4.4	4d)	35
4.4.5	4e)	35
4.4.6	4f)	35
4.4.7	4g)	36
4.4.8	4h)	36
4.4.9	4i)	36
4.4.10	4j)	36
4.4.11	4k)	36
4.4.12	4l)	36
4.4.13	4m)	37
4.4.14	4n)	37
4.5	5) Reproducing the Results with Scikit-Learn (15 points)	37
4.5.1	5a)	37

4.5.2	5b)	37
4.6	6) Bonus Question (15 points)	38
4.6.1	6a)	38
4.6.2	6b)	38
4.6.3	6c)	38
<b>5</b>	<b>Assignment 2</b>	<b>39</b>
	Instructions	39
5.1	1) Multiple Linear Regression (24 points)	40
5.1.1	1a)	40
5.1.2	1b)	40
5.1.3	1c)	40
5.1.4	1d)	40
5.1.5	1e)	40
5.1.6	1f)	41
5.1.7	1g)	41
5.1.8	1h)	41

# Preface

These are class notes for the course STAT303-2. This is not the course text-book. You are required to read the relevant sections of the book as mentioned on the course website.

The course notes are currently being written, and will continue to being developed as the course progresses (just like the course textbook last quarter). Please report any typos / mistakes / inconsistencies / issues with the class notes / class presentations in your comments [here](#). Thank you!

# **Part I**

## **Linear regression**

# 1 Simple Linear Regression

*Read section 3.1 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

## 1.1 Simple Linear Regression

```
import pandas as pd
import numpy as np
import statsmodels.formula.api as smf
import statsmodels.api as sm
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib.patches import Patch
from matplotlib.lines import Line2D
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
```

**Develop a simple linear regression model that predicts car price based on engine size.** Datasets to be used: *Car\_features\_train.csv*, *Car\_prices\_train.csv*

```
# We are reading training data ONLY at this point.
# Test data is already separated in another file
trainf = pd.read_csv('./Datasets/Car_features_train.csv') # Predictors
trainp = pd.read_csv('./Datasets/Car_prices_train.csv') # Response
train = pd.merge(trainf, trainp)
train.head()
```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

### 1.1.1 Training with `statsmodels`

Here, we will use the `statsmodels.formula.api` module of the `statsmodels` library. The use of “API” here doesn’t refer to a traditional external web API but rather an interface within the library for users to interact with and perform specific tasks. The `statsmodels.formula.api` module provides a formulaic interface to the `statsmodels` library. A formula is a compact way to specify statistical models using a formula language. This module allows users to define statistical models using formulas similar to those used in R.

So, in summary, the `statsmodels.formula.api` module provides a formulaic interface as part of the `statsmodels` library, allowing users to specify statistical models using a convenient and concise formula syntax.

```
# Let's create the model

# ols stands for Ordinary Least Squares - the name of the algorithm that optimizes Linear Regression

# data input needs the dataframe that has the predictor and the response
# formula input needs to:
#   # be a string
#   # have the following syntax: "response~predictor"

# Using engineSize to predict price
ols_object = smf.ols(formula = 'price~engineSize', data = train)

#Using the fit() function of the 'ols' class to fit the model, i.e., train the model
model = ols_object.fit()

#Printing model summary which contains among other things, the model coefficients
model.summary()
```



<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.390
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.390
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	3177.
<b>Date:</b>	Tue, 16 Jan 2024	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	16:46:33	<b>Log-Likelihood:</b>	-53949.
<b>No. Observations:</b>	4960	<b>AIC:</b>	1.079e+05
<b>Df Residuals:</b>	4958	<b>BIC:</b>	1.079e+05
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	-4122.0357	522.260	-7.893	0.000	-5145.896	-3098.176
<b>engineSize</b>	1.299e+04	230.450	56.361	0.000	1.25e+04	1.34e+04

<b>Omnibus:</b>	1271.986	<b>Durbin-Watson:</b>	0.517
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	6490.719
<b>Skew:</b>	1.137	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	8.122	<b>Cond. No.</b>	7.64

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

The model equation is:  $\hat{price} = -4122.0357 + 12990 * engineSize$

- R-squared is 39%. This is the proportion of variance in car price explained by `engineSize`.
- The coef of `engineSize` ( $\hat{\beta}_1$ ) is statistically significant ( $p$ -value = 0). There is a linear relationship between X and Y.
- The 95% CI of  $\hat{\beta}_1$  is [1.25e+04, 1.34e+04].
- PI is not shown here.

The coefficient of `engineSize` is 1.299e+04. - Unit change in `engineSize` increases the expected price by \$ 12,990. - An increase of 3 increases the price by \$  $(3 * 1.299e+04) = \$ 38,970$ .

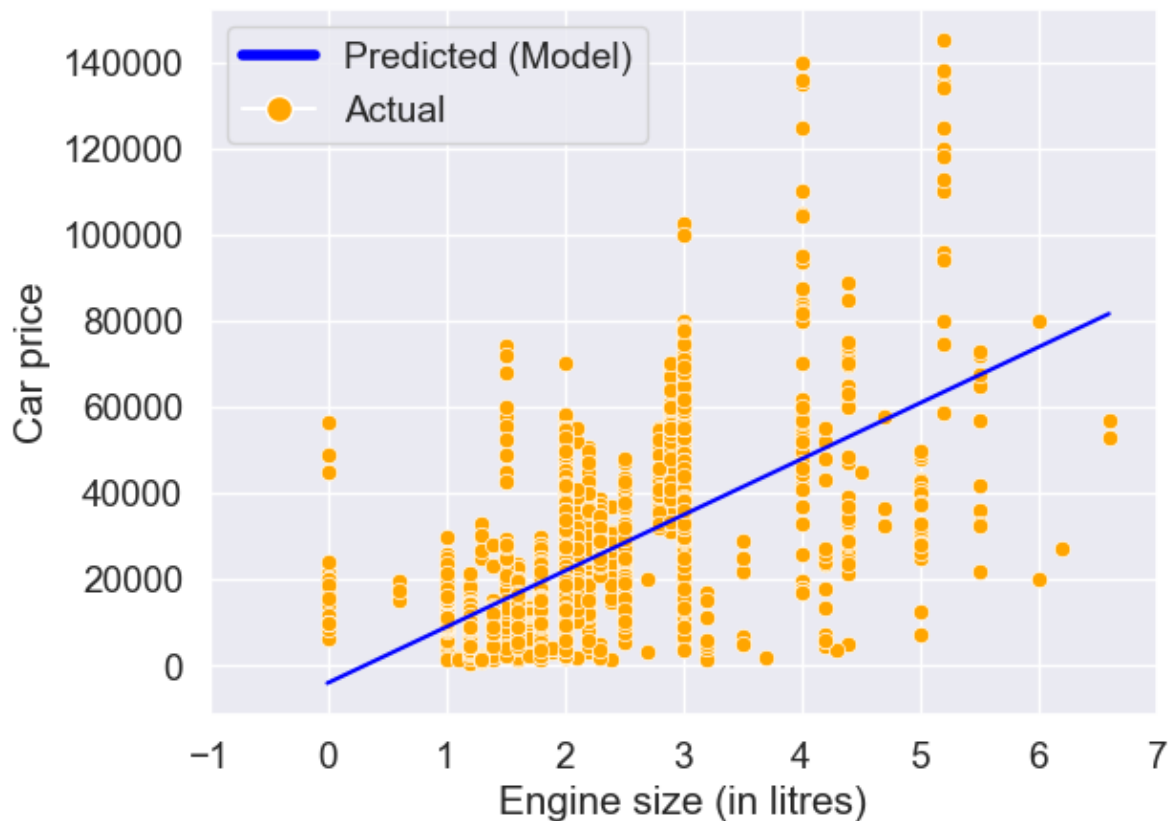
The coefficients can also be returned directly using the `params` attribute of the `model` object returned by the `fit()` method of the `ols` class:

```
model.params
```

```
Intercept    -4122.035744
engineSize    12988.281021
dtype: float64
```

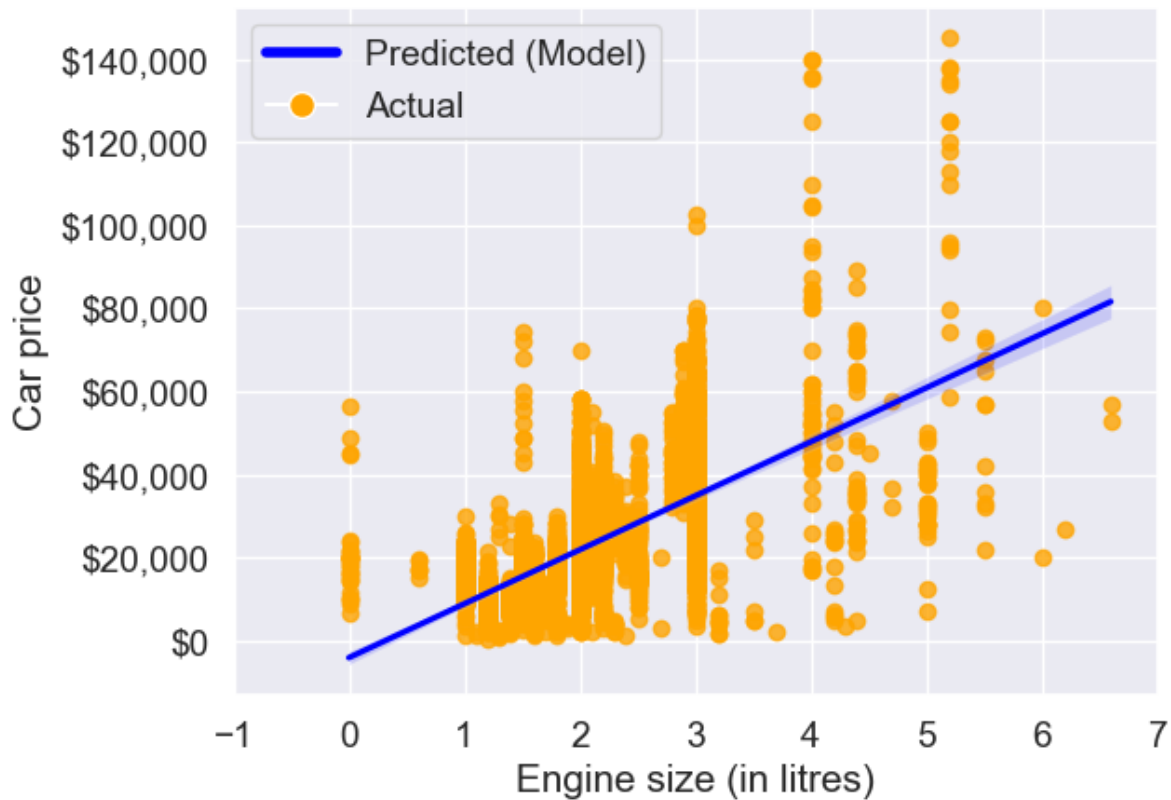
## Visualize the regression line

```
sns.set(font_scale=1.25)
ax = sns.scatterplot(x = train.engineSize, y = train.price,color = 'orange')
sns.lineplot(x = train.engineSize, y = model.fittedvalues,color = 'blue')
plt.xlim(-1,7)
plt.xlabel('Engine size (in litres)')
plt.ylabel('Car price')
legend_elements = [Line2D([0], [0], color='blue', lw=4, label='Predicted (Model)'),
                   Line2D([0], [0], marker='o', color='w', label='Actual',
                           markerfacecolor='orange', markersize=10)]
ax.legend(handles=legend_elements, loc='upper left');
```



Note that the above plot can be made directly using the seaborn function `regplot()`. The function `regplot()` fits a simple linear regression model with `y` as the response, and `x` as the predictor, and then plots the model over a scatterplot of the data.

```
ax = sns.regplot(x = 'engineSize', y = 'price', data = train, color = 'orange', line_kws={"co":
plt.xlim(-1,7)
plt.xlabel('Engine size (in litres)')
plt.ylabel('Car price')
ax.yaxis.set_major_formatter('${x:,.0f}')
ax.legend(handles=legend_elements, loc='upper left');
#Note that some of the engineSize values are 0. They are incorrect, and should ideally be im
```



The light shaded region around the blue line in the above plot is the confidence interval.

**Predict the car price for the cars in the test dataset.** Datasets to be used: *Car\_features\_test.csv*, *Car\_prices\_test.csv*

Now that the model has been trained, let us evaluate it on unseen data. Make sure that the columns names of the predictors are the same in train and test datasets.

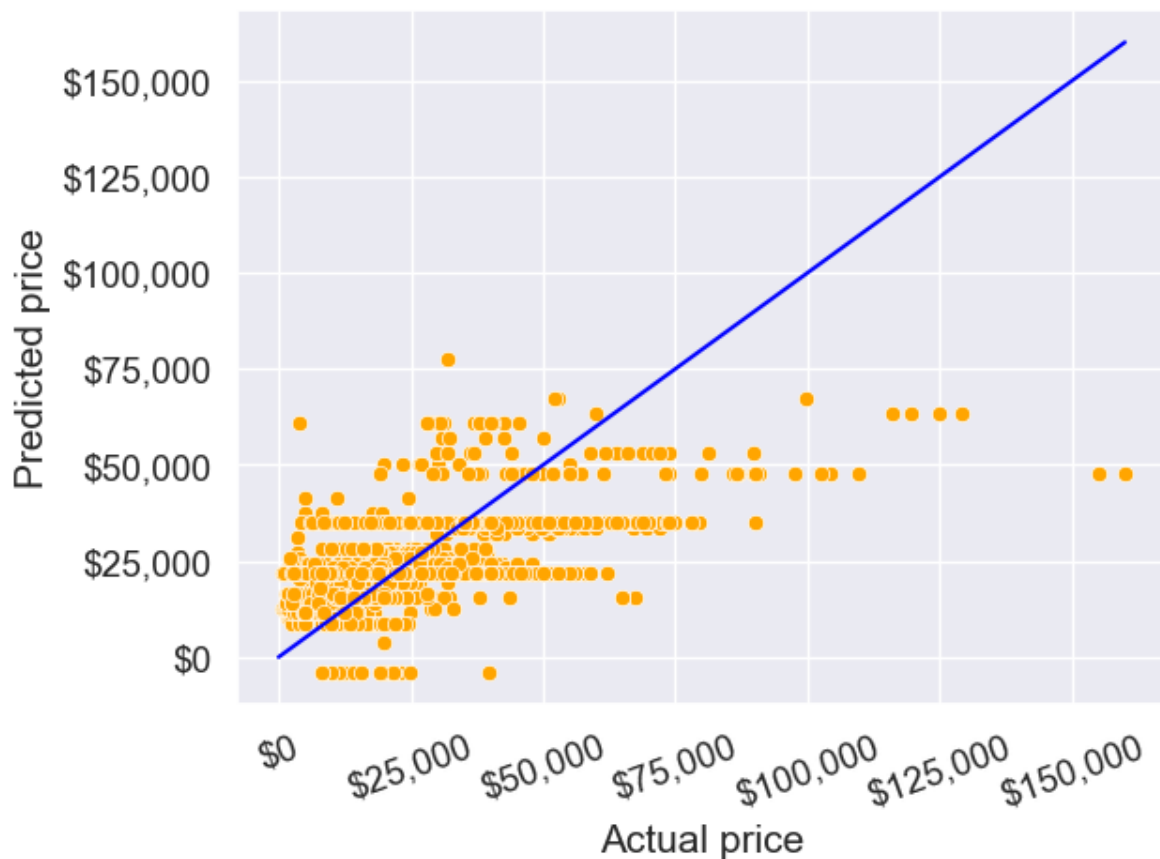
```
# Read the test data
testf = pd.read_csv('./Datasets/Car_features_test.csv') # Predictors
```

```
testp = pd.read_csv('./Datasets/Car_prices_test.csv') # Response
test = pd.merge(testf, testp)
```

```
#Using the predict() function associated with the 'model' object to make predictions of car p
pred_price = model.predict(testf)#Note that the predict() function finds the predictor 'engi
```

**Make a visualization that compares the predicted car prices with the actual car prices**

```
sns.scatterplot(x = testp.price, y = pred_price, color = 'orange')
#In case of a perfect prediction, all the points must lie on the line x = y.
ax = sns.lineplot(x = [0,testp.price.max()], y = [0,testp.price.max()],color='blue') #Plottin
plt.xlabel('Actual price')
plt.ylabel('Predicted price')
ax.yaxis.set_major_formatter('${x:,.0f}')
ax.xaxis.set_major_formatter('${x:,.0f}')
plt.xticks(rotation=20);
```



The prediction doesn't look too good. This is because we are just using one predictor - engine size. We can probably improve the model by adding more predictors when we learn multiple linear regression.

**What is the RMSE of the predicted car price on unseen data?**

```
np.sqrt(((testp.price - pred_price)**2).mean())
```

12995.106451548696

The root mean squared error in predicting car price is around \$13k.

**What is the residual standard error based on the training data?**

```
np.sqrt(model.mse_resid)
```

12810.109175214138

The residual standard error on the training data is close to the RMSE on the test data. This shows that the performance of the model on unknown data is comparable to its performance on known data. This implies that the model is not overfitting, which is good! In case we overfit a model on the training data, its performance on unknown data is likely to be worse than that on the training data.

**Find the confidence and prediction intervals of the predicted car price**

```
#Using the get_prediction() function associated with the 'model' object to get the intervals
intervals = model.get_prediction(testf)
```

```
#The function requires specifying alpha (probability of Type 1 error) instead of the confidence
intervals.summary_frame(alpha=0.05)
```

	mean	mean_se	mean_ci_lower	mean_ci_upper	obs_ci_lower	obs_ci_upper
0	34842.807319	271.666459	34310.220826	35375.393812	9723.677232	59961.937406
1	34842.807319	271.666459	34310.220826	35375.393812	9723.677232	59961.937406
2	34842.807319	271.666459	34310.220826	35375.393812	9723.677232	59961.937406
3	8866.245277	316.580850	8245.606701	9486.883853	-16254.905974	33987.396528
4	47831.088340	468.949360	46911.740050	48750.436631	22700.782946	72961.393735
...	...	...	...	...	...	...
2667	47831.088340	468.949360	46911.740050	48750.436631	22700.782946	72961.393735

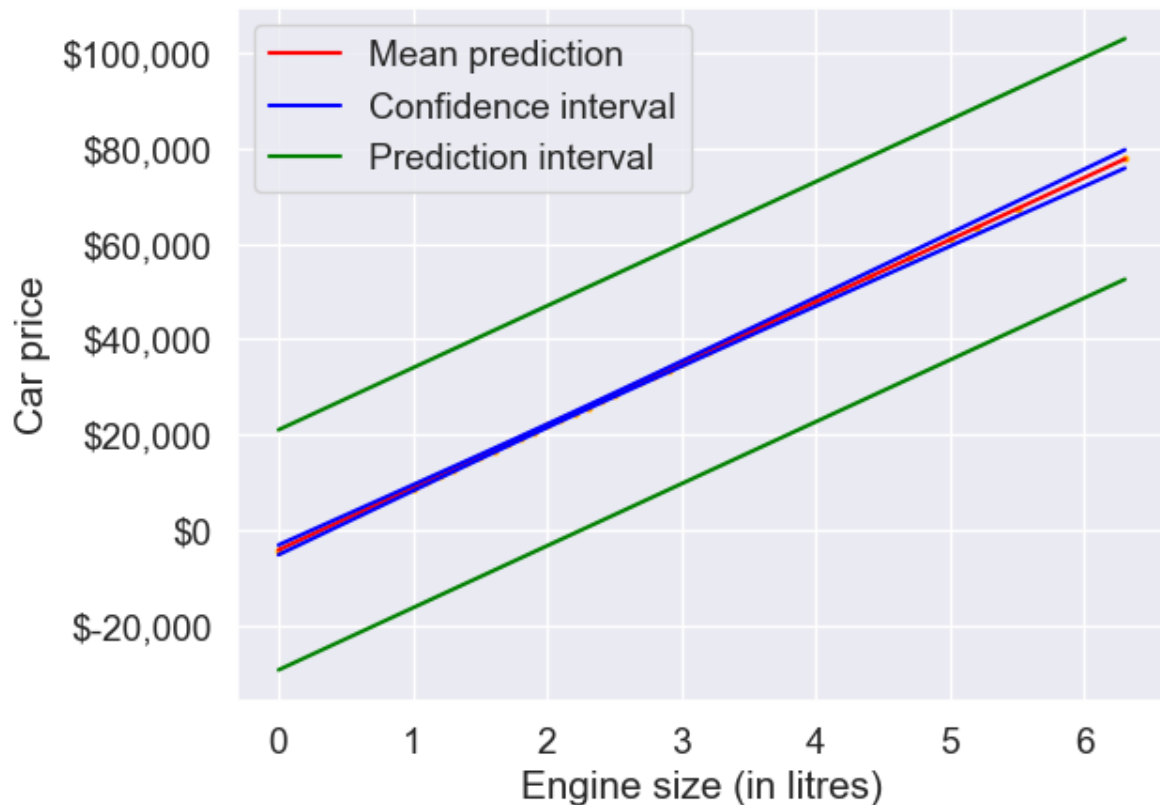
	mean	mean_se	mean_ci_lower	mean_ci_upper	obs_ci_lower	obs_ci_upper
2668	34842.807319	271.666459	34310.220826	35375.393812	9723.677232	59961.937406
2669	8866.245277	316.580850	8245.606701	9486.883853	-16254.905974	33987.396528
2670	21854.526298	184.135754	21493.538727	22215.513869	-3261.551421	46970.604017
2671	21854.526298	184.135754	21493.538727	22215.513869	-3261.551421	46970.604017

Show the regression line predicting car price based on engine size for test data.  
Also show the confidence and prediction intervals for the car price.

```
interval_table = intervals.summary_frame(alpha=0.05)

ax = sns.scatterplot(x = testf.engineSize, y = pred_price, color = 'orange', s = 10)
sns.lineplot(x = testf.engineSize, y = pred_price, color = 'red')
sns.lineplot(x = testf.engineSize, y = interval_table.mean_ci_lower, color = 'blue')
sns.lineplot(x = testf.engineSize, y = interval_table.mean_ci_upper, color = 'blue')
sns.lineplot(x = testf.engineSize, y = interval_table.obs_ci_lower, color = 'green')
sns.lineplot(x = testf.engineSize, y = interval_table.obs_ci_upper, color = 'green')

legend_elements = [Line2D([0], [0], color='red', label='Mean prediction'),
                    Line2D([0], [0], color='blue', label='Confidence interval'),
                    Line2D([0], [0], color='green', label='Prediction interval')]
ax.legend(handles=legend_elements, loc='upper left')
plt.xlabel('Engine size (in litres)')
plt.ylabel('Car price')
ax.yaxis.set_major_formatter('${x:,.0f}');
```



### 1.1.2 Training with `sklearn`

```
# Create the model as an object

model = LinearRegression() # No inputs, this will change for other models

# Train the model - separate the predictor(s) and the response for this!
X_train = train[['engineSize']]
y_train = train[['price']]

# Note that both are dfs, NOT series - necessary to avoid errors

model.fit(X_train, y_train)

# Check the slight syntax differences
# predictors and response separate
# We need to manually slice the predictor column(s) we want to include
```

```

# No need to assign to an output

# Return the parameters
print("Coefficient of engine size = ", model.coef_) # slope
print("Intercept = ", model.intercept_) # intercept

# No .summary() here! - impossible to do much inference; this is a shortcoming of sklearn

```

```

Coefficient of engine size =  [[12988.28102112]]
Intercept =  [-4122.03574424]

```

```

# Prediction

# Again, separate the predictor(s) and the response of interest
X_test = test[['engineSize']]
y_test = test[['price']].to_numpy() # Easier to handle with calculations as np array

y_pred = model.predict(X_test)

# Evaluate
model_rmse = np.sqrt(np.mean((y_pred - y_test)**2)) # RMSE
model_mae = np.mean(np.abs(y_pred - y_test)) # MAE

print('Test RMSE: ', model_rmse)

```

```

Test RMSE:  12995.106451548696

```

```

# Easier way to calculate metrics with sklearn tools

# Note that we have imported the functions 'mean_squared_error' and 'mean_absolute_error'
# from the sklearn.metrics module (check top of the code)

model_rmse = np.sqrt(mean_squared_error(y_test,y_pred))
model_mae = mean_absolute_error(y_test,y_pred)
print('Test RMSE: ', model_rmse)
print('Test MAE: ', model_mae)

```

```

Test RMSE:  12995.106451548696
Test MAE:   9411.325912951994

```



```

y_pred_train = model.predict(X_train)
print('Train R-squared:', r2_score(y_train, y_pred_train))
print('Test R-squared:', r2_score(y_test, y_pred))

```

Train R-squared: 0.39049842625794573

Test R-squared: 0.3869900378620146

**Note:** Why did we repeat the same task in two different libraries?

- `statsmodels` and `sklearn` have different advantages - we will use both for our purposes
  - `statsmodels` returns a lot of statistical output, which is very helpful for inference (coming up next) but it has a limited variety of models.
  - With `statsmodels`, you may have columns in your DataFrame in addition to predictors and response, while with `sklearn` you need to make separate objects consisting of only the predictors and the response.
  - `sklearn` includes many models (Lasso and Ridge this quarter, many others next quarter) and helpful tools/functions (like metrics) that `statsmodels` does not but it does not have any inference tools.

### 1.1.3 Training with `statsmodels.api`

Earlier we had used the `statsmodels.formula.api` module, where we had to put the regression model as a formula. We can also use the `statsmodels.api` module to develop a regression model. The syntax of training a model with the `OLS()` function in this module is similar to that of `sklearn`'s `LinearRegression()` function. However, the order in which the predictors and response are specified is different. The formula-style syntax of the `statsmodels.formula.api` module is generally preferred. However, depending on the situation, the `OLS()` syntax of `statsmodels.api` may be preferred.

Note that you will manually need to add the predictor (*a column of ones*) corresponding to the intercept to train the model with this method.

```

# Create the model as an object

# Train the model - separate the predictor(s) and the response for this!
X_train = train[['engineSize']]
y_train = train[['price']]

X_train_with_intercept = np.concatenate((np.ones(X_train.shape[0]).reshape(-1,1), X_train), a

model = sm.OLS(y_train, X_train_with_intercept).fit()

```

```
# Return the parameters
print(model.params)
```

```
const    -4122.035744
x1        12988.281021
dtype: float64
```

The model summary and all other attributes and methods of the `model` object are the same as that with the object created using the `statsmodels.formula.api` module.

```
model.summary()
```

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.390
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.390
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	3177.
<b>Date:</b>	Mon, 08 Jan 2024	<b>Prob (F-statistic):</b>	0.00
<b>Time:</b>	11:17:55	<b>Log-Likelihood:</b>	-53949.
<b>No. Observations:</b>	4960	<b>AIC:</b>	1.079e+05
<b>Df Residuals:</b>	4958	<b>BIC:</b>	1.079e+05
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P>  t	[0.025	0.975]
<b>const</b>	-4122.0357	522.260	-7.893	0.000	-5145.896	-3098.176
<b>x1</b>	1.299e+04	230.450	56.361	0.000	1.25e+04	1.34e+04

<b>Omnibus:</b>	1271.986	<b>Durbin-Watson:</b>	0.517
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	6490.719
<b>Skew:</b>	1.137	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	8.122	<b>Cond. No.</b>	7.64

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

## 2 Multiple Linear Regression

*Read section 3.2 of the book before using these notes.*

*Note that in this course, lecture notes are not sufficient, you must read the book for better understanding. Lecture notes are just implementing the concepts of the book on a dataset, but not explaining the concepts elaborately.*

### 2.1 Multiple Linear Regression

```
# importing libraries
import pandas as pd
import numpy as np
import statsmodels.formula.api as smf
import seaborn as sns
import matplotlib.pyplot as plt
```

**Develop a multiple linear regression model that predicts car price based on engine size, year, mileage, and mpg.** Datasets to be used: *Car\_features\_train.csv*, *Car\_prices\_train.csv*

```
# Reading datasets
trainf = pd.read_csv('./Datasets/Car_features_train.csv')
trainp = pd.read_csv('./Datasets/Car_prices_train.csv')
train = pd.merge(trainf, trainp)
train.head()
```

	carID	brand	model	year	transmission	mileage	fuelType	tax	mpg	engineSize	price
0	18473	bmw	6 Series	2020	Semi-Auto	11	Diesel	145	53.3282	3.0	37980
1	15064	bmw	6 Series	2019	Semi-Auto	10813	Diesel	145	53.0430	3.0	33980
2	18268	bmw	6 Series	2020	Semi-Auto	6	Diesel	145	53.4379	3.0	36850
3	18480	bmw	6 Series	2017	Semi-Auto	18895	Diesel	145	51.5140	3.0	25998
4	18492	bmw	6 Series	2015	Automatic	62953	Diesel	160	51.4903	3.0	18990

```
#Using the ols function to create an ols object. 'ols' stands for 'Ordinary least squares'
ols_object = smf.ols(formula = 'price~year+mileage+mpg+engineSize', data = train)
model = ols_object.fit()
model.summary()
```

Table 2.2: OLS Regression Results

Dep. Variable:	price	R-squared:	0.660
Model:	OLS	Adj. R-squared:	0.660
Method:	Least Squares	F-statistic:	2410.
Date:	Tue, 27 Dec 2022	Prob (F-statistic):	0.00
Time:	01:07:25	Log-Likelihood:	-52497.
No. Observations:	4960	AIC:	1.050e+05
Df Residuals:	4955	BIC:	1.050e+05
Df Model:	4		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	-3.661e+06	1.49e+05	-24.593	0.000	-3.95e+06	-3.37e+06
year	1817.7366	73.751	24.647	0.000	1673.151	1962.322
mileage	-0.1474	0.009	-16.817	0.000	-0.165	-0.130
mpg	-79.3126	9.338	-8.493	0.000	-97.620	-61.006
engineSize	1.218e+04	189.969	64.107	0.000	1.18e+04	1.26e+04

Omnibus:	2450.973	Durbin-Watson:	0.541
Prob(Omnibus):	0.000	Jarque-Bera (JB):	31060.548
Skew:	2.045	Prob(JB):	0.00
Kurtosis:	14.557	Cond. No.	3.83e+07

The model equation is: estimated car price =  $-3.661e6 + 1818 * \text{year} - 0.15 * \text{mileage} - 79.31 * \text{mpg} + 12180 * \text{engineSize}$

**Predict the car price for the cars in the test dataset.** Datasets to be used: *Car\_features\_test.csv*, *Car\_prices\_test.csv*

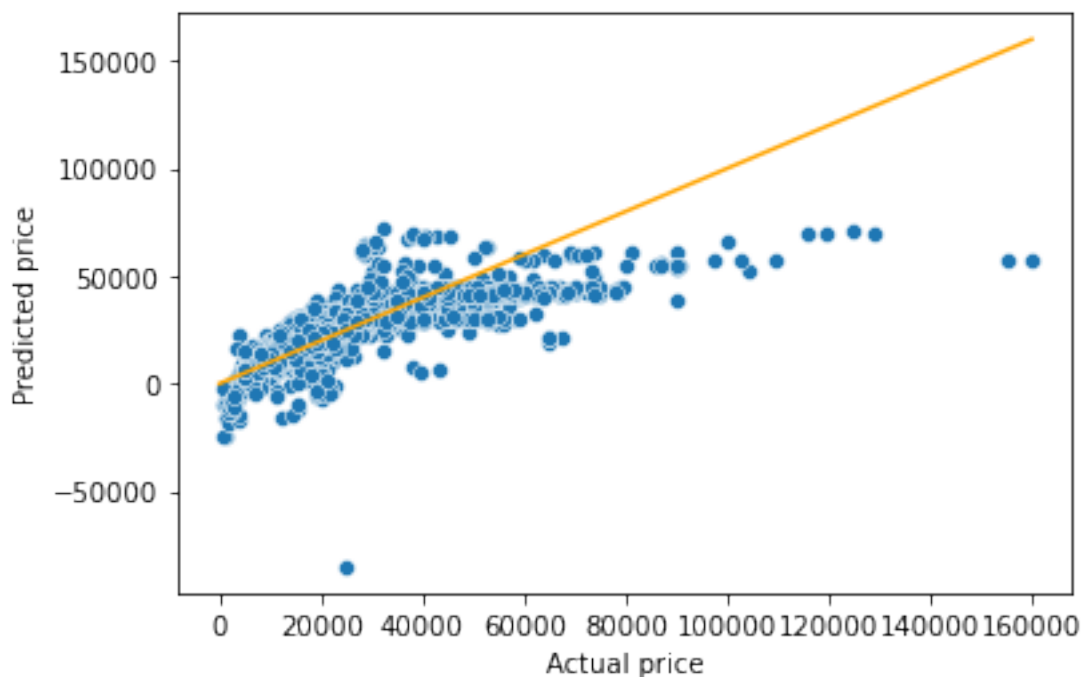
```
testf = pd.read_csv('./Datasets/Car_features_test.csv')
testp = pd.read_csv('./Datasets/Car_prices_test.csv')
```

```
#Using the predict() function associated with the 'model' object to make predictions of car price
pred_price = model.predict(testf)#Note that the predict() function finds the predictor 'engine'
```

**Make a visualization that compares the predicted car prices with the actual car prices**

```
sns.scatterplot(x = testp.price, y = pred_price)
#In case of a perfect prediction, all the points must lie on the line x = y.
sns.lineplot(x = [0,testp.price.max()], y = [0,testp.price.max()],color='orange') #Plotting the line x = y
plt.xlabel('Actual price')
plt.ylabel('Predicted price')
```

```
Text(0, 0.5, 'Predicted price')
```



The prediction looks better as compared to the one with simple linear regression. This is because we have four predictors to help explain the variation in car price, instead of just one in the case of simple linear regression. Also, all the predictors have a significant relationship with price as evident from their p-values. Thus, all four of them are contributing in explaining the variation. Note the higher values of  $R^2$  as compared to the one in the case of simple linear regression.

**What is the RMSE of the predicted car price?**

```
np.sqrt(((testp.price - pred_price)**2).mean())
```

9956.82497993548

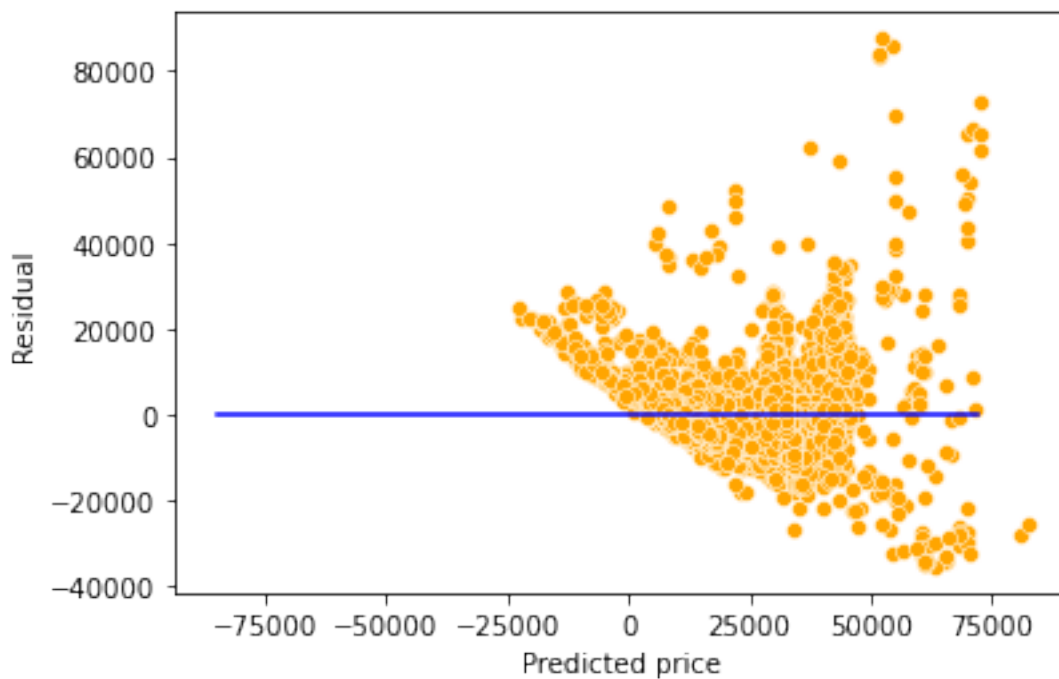
What is the residual standard error based on the training data?

```
np.sqrt(model.mse_resid)
```

9563.74782917604

```
sns.scatterplot(x = model.fittedvalues, y=model.resid,color = 'orange')
sns.lineplot(x = [pred_price.min(),pred_price.max()],y = [0,0],color = 'blue')
plt.xlabel('Predicted price')
plt.ylabel('Residual')
```

```
Text(0, 0.5, 'Residual')
```



Will the explained variation (R-squared) in car price always increase if we add a variable?

Should we keep on adding variables as long as the explained variation (R-squared) is increasing?

```
#Using the ols function to create an ols object. 'ols' stands for 'Ordinary least squares'
np.random.seed(1)
train['rand_col'] = np.random.rand(train.shape[0])
ols_object = smf.ols(formula = 'price~year+mileage+mpg+engineSize+rand_col', data = train)
model = ols_object.fit()
model.summary()
```

Table 2.5: OLS Regression Results

Dep. Variable:	price	R-squared:	0.661
Model:	OLS	Adj. R-squared:	0.660
Method:	Least Squares	F-statistic:	1928.
Date:	Tue, 27 Dec 2022	Prob (F-statistic):	0.00
Time:	01:07:38	Log-Likelihood:	-52497.
No. Observations:	4960	AIC:	1.050e+05
Df Residuals:	4954	BIC:	1.050e+05
Df Model:	5		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	-3.662e+06	1.49e+05	-24.600	0.000	-3.95e+06	-3.37e+06
year	1818.1672	73.753	24.652	0.000	1673.578	1962.756
mileage	-0.1474	0.009	-16.809	0.000	-0.165	-0.130
mpg	-79.2837	9.338	-8.490	0.000	-97.591	-60.976
engineSize	1.218e+04	189.972	64.109	0.000	1.18e+04	1.26e+04
rand_col	451.1226	471.897	0.956	0.339	-474.004	1376.249

Omnibus:	2451.728	Durbin-Watson:	0.541
Prob(Omnibus):	0.000	Jarque-Bera (JB):	31040.331
Skew:	2.046	Prob(JB):	0.00
Kurtosis:	14.552	Cond. No.	3.83e+07

Adding a variable with random values to the model (*rand\_col*) increased the explained variation (R-squared). This is because the model has one more parameter to tune to reduce the residual squared error (RSS). However, the p-value of *rand\_col* suggests that its coefficient is zero. Thus, using the model with *rand\_col* may give poorer performance on unknown data, as compared to the model without *rand\_col*. This implies that it is not a good idea to blindly add variables in the model to increase R-squared.

## **Part II**

# **Assignments**



## 3 Assignment 1 (Section 20)

### Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Do not write your name on the assignment.
3. Write your code in the **Code cells** and your answers in the **Markdown cells** of the Jupyter notebook. Ensure that the solution is written neatly enough to for the graders to understand and follow.
4. Use [Quarto](#) to render the **.ipynb** file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
5. The assignment is worth 100 points, and is due on **Wednesday, 24th January 2024 at 11:59 pm**.
6. There is a **bonus** question worth 15 points.
7. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (**1 point**). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file.*
  - No name can be written on the assignment, nor can there be any indicator of the student's identity—e.g. printouts of the working directory should not be included in the final submission. (**1 point**)
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (**1 point**)
  - Final answers to each question are written in the Markdown cells. (**1 point**)
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text. (**1 point**)

8. The maximum possible score in the assignment is  $100 + 15$  (bonus question)  $+ 5$  (proper formatting) = 120 out of 100. There is no partial credit for some parts of the bonus question.

### **3.1 1) Case Studies: Regression vs Classification and Prediction vs Inference (16 points)**

#### **3.1.1 1a)**

For each case below, explain (1) whether it is a classification or a regression problem and (2) whether the main purpose is prediction or inference. **You need justify your answers for credit.**

#### **3.1.2 1b)**

You work for a company that is interested in conducting a marketing campaign. The goal of your project is to identify individuals who are likely to respond positively to a marketing campaign, based on observations of demographic variables (such as age, gender, income etc.) measured on each individual. **(2+2 points)**

#### **3.1.3 1c)**

For the same company, now you are working on a different project. This one is focused on understanding the impact of advertisements in different media types on the company sales. For example, you are interested in the following question: *'How large of an increase in sales is associated with a given increase in radio and TV advertising?'* **(2+2 points)**

#### **3.1.4 1d)**

A company is selling furniture and they are interested in the finding the association between demographic characteristics of customers (such as age, gender, income etc.) and if they would purchase a particular company product. **(2+2 points)**

### 3.1.5 1e)

We are interested in forecasting the % change in the USD/Euro exchange rate using the weekly changes in the stock markets of a number of countries. We collect weekly data for all of 2023. For each week, we record the % change in the USD/Euro, the % change in the US market, the % change in the British market, and the % change in the German market. **(2+2 points)**

## 3.2 2) Examples for Different Regression Metrics: RMSE vs MAE (8 points)

### 3.2.1 2a)

Describe a regression problem, where it will be more proper to evaluate the model performance using the root mean squared error (RMSE) metric as compared to the mean absolute error (MAE) metric. **You need to justify your answer for credit. (4 points)**

**Note:** You are not allowed to use the datasets and examples covered in the lectures.

### 3.2.2 2b)

Describe a regression problem, where it will be more proper to evaluate the model performance using the mean absolute error (MAE) metric as compared to the root mean squared error (RMSE) metric. **You need to justify your answer for credit. (4 points)**

**Note:** You are not allowed to use the datasets and examples covered in the lectures.

## 3.3 3) Modeling the Petrol Consumption in U.S. States (61 points)

Read `petrol_consumption_train.csv`. Assume that each observation is a U.S. state. For each observation, the data has the following variables as its five columns:

`Petrol_tax`: Petrol tax (cents per gallon)

`Per_capita_income`: Average income (dollars)

`Paved_highways`: Paved Highways (miles)

`Prop_license`: Proportion of population with driver's licenses

`Petrol_consumption`: Consumption of petrol (millions of gallons)

### 3.3.1 3a)

Create a pairwise plot of all the variables in the dataset. **(1 point)** Print the correlation matrix of all the variables as well. **(1 point)** Which variable has the highest linear correlation with `Petrol_consumption`? **(2 points)**

**Note:** Remember that a pairwise plot is a visualization tool that you can find in the seaborn library.

### 3.3.2 3b)

Fit a simple linear regression model to predict `Petrol_consumption` using the column you found in **part a** as the only predictor. Print the model summary. **(4 points)**

### 3.3.3 3c)

What is the increase in petrol consumption for an increase of 0.05 in the predictor? **(4 points)**

### 3.3.4 3d)

Does petrol consumption have a statistically significant relationship with the predictor? **You need to justify your answer for credit. (4 points)**

### 3.3.5 3e)

How much of the variation in petrol consumption can be explained by its linear relationship with the predictor? **(3 points)**

### 3.3.6 3f)

Predict the petrol consumption for a state in which 50% of the population has a driver's license. **(3 points)** What are the confidence interval **(3 points)** and the prediction interval **(3 points)** for your prediction? Which interval is wider? **(1 points)** Why? **(2 points)**

### 3.3.7 3g)

Predict the petrol consumption for a state in which 10% of the population has a driver's license. **(3 points)** Are you getting a reasonable outcome? **(1 point)** Why or why not? **(2 points)**

### 3.3.8 3h)

What is the residual standard error of the model? **(3 points)**

### 3.3.9 3i)

Using the trained model, predict the petrol consumption of the observations in `petrol_consumption_test.csv`. **(2 points)** and find the RMSE. **(2 points)** What is the unit of this RMSE value? **(1 point)**

### 3.3.10 3j)

Based on the answers to **part g** and **part h**, do you think the model is overfitting? **You need to justify your answer for credit. (4 points)**

### 3.3.11 3k)

Make a scatterplot of `Petrol_consumption` vs. the predictor using `petrol_consumption_test.csv`. **(1 point)** Over the scatterplot, plot the regression line **(2 points)**, the prediction interval **(2 points)**, and the confidence interval. **(2 points)**

Make sure that regression line, prediction interval lines, and confidence interval lines have different colors. **(1 point)** Display a legend that correctly labels the lines as well. **(1 point)** Note that you need two lines of the same color to plot an interval.

### 3.3.12 3l)

Find the correlation between `Petrol_consumption` and the rest of the variables in `petrol_consumption_train.csv`. Which column would have the lowest R-squared value when used as the predictor for a Simple Linear Regression model to predict `Petrol_consumption`? Note that you can directly answer this question from the correlation values and do not need to develop any more linear regression models. **(3 points)**

## 3.4 4) Reproducing the Results with Scikit-Learn (15 points)

### 3.4.1 4a)

Using the same datasets, same response and the same predictor as **Question 3**, reproduce the following outputs with scikit-learn:

- Model RMSE for test data (**3 points**)
- R-squared value of the model (**3 points**)
- Residual standard error of the model (**3 points**)

**Note that you are only allowed to use scikit-learn, pandas, and numpy tools for this question. Any other libraries will not receive any credit.**

### 3.4.2 4b)

Which of the model outputs from **Question 3** cannot be reproduced using scikit-learn? Give two answers. (**2+2 points**) What does this tell about scikit-learn? (**2 points**)

## 3.5 5) Bonus Question (15 points)

Please note that the bonus question requires you to look more into the usage of the tools we covered in class and it will be necessary to do your own research. We strongly suggest attempting it after you are done with the rest of the assignment.

### 3.5.1 5a)

Fit a simple linear regression model to predict `Petrol_consumption` based on the predictor in **Question 3**, but **without an intercept term**. (**5 points - no partial credit**)

Without an intercept means that the equation becomes  $Y = \beta_1 X$ . The intercept term,  $\beta_0$ , becomes 0.

**Note: You must answer this part correctly to qualify for the bonus points in the following parts.**

### 3.5.2 5b)

Predict the petrol consumption for the observations in `petrol_consumption_test.csv` using the model without an intercept and find the RMSE. (**1+2 points**) Then, print the summary and find the R-squared. (**2 points**)

### 3.5.3 5c)

The RMSE for the models with and without the intercept are similar, which indicates that both models are almost equally good. However, the R-squared for the model without intercept is much higher than the R-squared for the model with the intercept. Why? Justify your answer. **(5 points - no partial credit)**

## 4 Assignment 1 (Sections 21 & 22)

### Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Do not write your name on the assignment.
3. Write your code in the **Code cells** and your answers in the **Markdown cells** of the Jupyter notebook. Ensure that the solution is written neatly enough to for the graders to understand and follow.
4. Use [Quarto](#) to render the **.ipynb** file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
5. The assignment is worth 100 points, and is due on **Wednesday, 24th January 2024 at 11:59 pm**.
6. There is a **bonus** question worth 15 points.
7. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (**1 point**). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file.*
  - No name can be written on the assignment, nor can there be any indicator of the student's identity—e.g. printouts of the working directory should not be included in the final submission. (**1 point**)
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (**1 point**)
  - Final answers to each question are written in the Markdown cells. (**1 point**)
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text. (**1 point**)



8. The maximum possible score in the assignment is  $100 + 15$  (bonus question)  $+ 5$  (proper formatting) = 120 out of 100. There is no partial credit for some parts of the bonus question.

## 4.1 1) Case Studies: Regression vs Classification and Prediction vs Inference (16 points)

For each case below, explain (1) whether it is a classification or a regression problem and (2) whether the main purpose is prediction or inference. **You need justify your answers for credit.**

### 4.1.1 1a)

You work for a company that is interested in conducting a marketing campaign. The goal of your project is to identify individuals who are likely to respond positively to a marketing campaign, based on observations of demographic variables (such as age, gender, income etc.) measured on each individual. **(2+2 points)**

### 4.1.2 1b)

For the same company, now you are working on a different project. This one is focused on understanding the impact of advertisements in different media types on the company sales. For example, you are interested in the following question: *‘How large of an increase in sales is associated with a given increase in radio and TV advertising?’* **(2+2 points)**

### 4.1.3 1c)

A company is selling furniture and they are interested in the finding the association between demographic characteristics of customers (such as age, gender, income etc.) and if they would purchase a particular company product. **(2+2 points)**

### 4.1.4 1d)

We are interested in forecasting the % change in the USD/Euro exchange rate using the weekly changes in the stock markets of a number of countries. We collect weekly data for all of 2023. For each week, we record the % change in the USD/Euro, the % change in the US market, the % change in the British market, and the % change in the German market. **(2+2 points)**

## 4.2 2) Examples for Different Regression Metrics: RMSE vs MAE (8 points)

### 4.2.1 2a)

Describe a regression problem, where it will be more proper to evaluate the model performance using the root mean squared error (RMSE) metric as compared to the mean absolute error (MAE) metric. **You need to justify your answer for credit. (4 points)**

**Note:** You are not allowed to use the datasets and examples covered in the lectures.

### 4.2.2 2b)

Describe a regression problem, where it will be more proper to evaluate the model performance using the mean absolute error (MAE) metric as compared to the root mean squared error (RMSE) metric. **You need to justify your answer for credit. (4 points)**

**Note:** You are not allowed to use the datasets and examples covered in the lectures.

## 4.3 3) Simple Linear Regression: Formulation (3 points)

When asked to state the simple linear regression model, a student wrote it as follows:  $E(Y_i) = \beta_0 + \beta_1 X_i + \epsilon_i$ . Is this correct (1 point)? Justify your answer (2 points).

## 4.4 4) Modeling the Petrol Consumption in U.S. States (58 points)

Read `petrol_consumption_train.csv`. Assume that each observation is a U.S. state. For each observation, the data has the following variables as its five columns:

`Petrol_tax`: Petrol tax (cents per gallon)

`Per_capita_income`: Average income (dollars)

`Paved_highways`: Paved Highways (miles)

`Prop_license`: Proportion of population with driver's licenses

`Petrol_consumption`: Consumption of petrol (millions of gallons)

#### 4.4.1 4a)

Create a pairwise plot of all the variables in the dataset. **(1 point)** Print the correlation matrix of all the variables as well. **(1 point)** Which variable has the highest linear correlation with `Petrol_consumption`? **(1 point)**

**Note:** Remember that a pairwise plot is a visualization tool that you can find in the seaborn library.

#### 4.4.2 4b)

Fit a simple linear regression model to predict `Petrol_consumption` using the column you found in **part a** as the only predictor. Print the model summary. **(3 points)**

#### 4.4.3 4c)

When asked for a point estimate of the expected petrol consumption for a state where the proportion of population with driver's license is 54.4%, a person gave the estimate 488 million gallons because that is the mean value of `Petrol_consumption` for the two observations of `Prop_license` = 0.544 pieces in the dataset. Is there an issue with this approach? Explain. **(2 points)** If there is an issue, then suggest a better approach and use it to estimate the expected petrol consumption for a state where the proportion of population with driver's license is 54.4%. **(2 points)**

#### 4.4.4 4d)

What is the increase in petrol consumption for an increase of 0.05 in the predictor? **(3 points)**

#### 4.4.5 4e)

Does petrol consumption have a statistically significant relationship with the predictor? **You need to justify your answer for credit. (3 points)**

#### 4.4.6 4f)

How much of the variation in petrol consumption can be explained by its linear relationship with the predictor? **(2 points)**

#### 4.4.7 4g)

Predict the petrol consumption for a state in which 50% of the population has a driver's license. **(2 points)** What are the confidence interval **(2 points)** and the prediction interval **(2 points)** for your prediction? Which interval is wider? **(1 points)** Why? **(2 points)**

#### 4.4.8 4h)

Predict the petrol consumption for a state in which 10% of the population has a driver's license. **(3 points)** Are you getting a reasonable outcome? **(1 point)** Why or why not? **(2 points)**

#### 4.4.9 4i)

What is the residual standard error of the model? **(3 points)**

#### 4.4.10 4j)

Using the trained model, predict the petrol consumption of the observations in `petrol_consumption_test.csv`. **(2 points)** and find the RMSE. **(2 points)** What is the unit of this RMSE value? **(1 point)**

#### 4.4.11 4k)

Based on the answers to **part i** and **part j**, do you think the model is overfitting? **You need to justify your answer for credit. (3 points)**

#### 4.4.12 4l)

Make a scatterplot of `Petrol_consumption` vs. the predictor using `petrol_consumption_test.csv`. **(1 point)** Over the scatterplot, plot the regression line **(1 point)**, the prediction interval **(2 points)**, and the confidence interval. **(2 points)**

Make sure that regression line, prediction interval lines, and confidence interval lines have different colors. **(1 point)** Display a legend that correctly labels the lines as well. **(1 point)** Note that you need two lines of the same color to plot an interval.

#### 4.4.13 4m)

The dataset consists of 40 US States. If you combine this data with the data of the remaining 10 US States, are you likely to obtain narrower confidence and prediction intervals in the plot developed in the previous question, for the same level of confidence? Justify your answer. (2 points).

If yes, then can you guarantee that the width of these intervals will reduce? Justify your answer. If no, then can you guarantee that the width of these intervals will not reduce? Justify your answer. (2 points)

#### 4.4.14 4n)

Find the correlation between `Petrol_consumption` and the rest of the variables in `petrol_consumption_train.csv`. Which column would have the lowest R-squared value when used as the predictor for a Simple Linear Regression model to predict `Petrol_consumption`? Note that you can directly answer this question from the correlation values and do not need to develop any more linear regression models. (2 points)

### 4.5 5) Reproducing the Results with Scikit-Learn (15 points)

#### 4.5.1 5a)

Using the same datasets, same response and the same predictor as **Question 4**, reproduce the following outputs with scikit-learn:

- Model RMSE for test data (3 points)
- R-squared value of the model (3 points)
- Residual standard error of the model (3 points)

Note that you are only allowed to use scikit-learn, pandas, and numpy tools for this question. Any other libraries will not receive any credit.

#### 4.5.2 5b)

Which of the model outputs from **Question 4** cannot be reproduced using scikit-learn? Give two answers. (2+2 points) What does this tell about scikit-learn? (2 points)

## 4.6 6) Bonus Question (15 points)

Please note that the bonus question requires you to look more into the usage of the tools we covered in class and it will be necessary to do your own research. We strongly suggest attempting it after you are done with the rest of the assignment.

### 4.6.1 6a)

Fit a simple linear regression model to predict `Petrol_consumption` based on the predictor in **Question 4**, but **without an intercept term**. (5 points - no partial credit)

Without an intercept means that the equation becomes  $Y = \beta_1 X$ . The intercept term,  $\beta_0$ , becomes 0.

**Note: You must answer this part correctly to qualify for the bonus points in the following parts.**

### 4.6.2 6b)

Predict the petrol consumption for the observations in `petrol_consumption_test.csv` using the model without an intercept and find the RMSE. (1+2 points) Then, print the summary and find the R-squared. (2 points)

### 4.6.3 6c)

The RMSE for the models with and without the intercept are similar, which indicates that both models are almost equally good. However, the R-squared for the model without intercept is much higher than the R-squared for the model with the intercept. Why? Justify your answer. (5 points - no partial credit)

## 5 Assignment 2

### Instructions

1. You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and not as a group activity.
2. Write your code in the **Code cells** and your answers in the **Markdown cells** of the Jupyter notebook. Ensure that the solution is written neatly enough to for the graders to understand and follow.
3. Use [Quarto](#) to render the **.ipynb** file as HTML. You will need to open the command prompt, navigate to the directory containing the file, and use the command: `quarto render filename.ipynb --to html`. Submit the HTML file.
4. The assignment is worth 100 points, and is due on **Sunday, 4th February 2024 at 11:59 pm**.
5. **Five points are properly formatting the assignment.** The breakdown is as follows:
  - Must be an HTML file rendered using Quarto (**1 point**). *If you have a Quarto issue, you must mention the issue & quote the error you get when rendering using Quarto in the comments section of Canvas, and submit the ipynb file.*
  - No name can be written on the assignment, nor can there be any indicator of the student's identity—e.g. printouts of the working directory should not be included in the final submission. (**1 point**)
  - There aren't excessively long outputs of extraneous information (e.g. no printouts of entire data frames without good reason, there aren't long printouts of which iteration a loop is on, there aren't long sections of commented-out code, etc.) (**1 point**)
  - Final answers to each question are written in the Markdown cells. (**1 point**)
  - There is no piece of unnecessary / redundant code, and no unnecessary / redundant text. (**1 point**)
6. The maximum possible score in the assignment is  $105 + 5$  (proper formatting) = 110 out of 100.

## 5.1 1) Multiple Linear Regression (24 points)

A study was conducted on 97 male patients with prostate cancer who were due to receive a radical prostatectomy (complete removal of the prostate). The **prostate.csv** file contains data on 9 measurements taken from these 97 patients. Each row (observation) represents a patient and each column (variable) represents a measurement. The description of variables can be found here: <https://rafalab.github.io/pages/649/prostate.html>

### 5.1.1 1a)

Fit a linear regression model with **lpsa** as the response and all the other variables as the predictors. Print its summary. (2 points) Write down the **optimal equation** that predicts **lpsa** using the predictors. (2 points)

### 5.1.2 1b)

Is the **overall regression** statistically significant? In other words, is there a statistically significant relationship between the response and at least one predictor? **You need to justify your answer for credit.** (2 points)

### 5.1.3 1c)

What does the optimal coefficient of **svi** tell us as a numeric output? Make sure you include the predictor, (**svi**) the response (**lpsa**) and the other predictors in your answer. (2 points)

### 5.1.4 1d)

Check the *p*-values of **gleason** and **age**. Are these predictors statistically significant? **You need to justify your answer for credit.** (2 points)

### 5.1.5 1e)

Check the 95% Confidence Interval of **age**. How can you relate it to its *p*-value and statistical significance, which you found in the previous part? (2 points)



### 5.1.6 1f)

This question requires some thinking, and bringing your 303-1 and 303-2 knowledge together.

Fit a **simple** linear regression model on `lpsa` against `gleason` and check the  $p$ -value of `gleason` using the summary. **(2 point)** Did the statistical significance of `gleason` change in the absence of other predictors? **(1 point)** Why or why not? **(3 points)**

**Hints:**

- 1) You need to compare this model with the Multiple Linear Regression model you created above.
- 2) Printing a correlation matrix of all the predictors should be useful.

### 5.1.7 1g)

Predict the `lpsa` of a **65 year old** man with `lcavol` = 1.35, `lweight` = 3.65, `lbph` = 0.1, `svi` = 0.22, `lcp` = -0.18, `gleason` = 6.75, and `pgg45` = 25. Find the 95% confidence and prediction intervals as well. **(2 points)**

### 5.1.8 1h)

In the Multiple Linear Regression model with all the predictors, you should see a total of five predictors that appear to be statistically insignificant. Why is it not a good idea to directly conclude that all of them are statistically insignificant? **(2 points)** Implement the additional test that concludes the statistical insignificance of all five predictors. **(2 points)**

**Hint:** `f_test()` method