

## Muhammad Salman Khan - 25373

### Machine Learning II - Assignment 2

#### Objective

The objective of this python notebook is to experiment on a regression problem. First, we will take a dataset and learn machine learning models on it and then we will find out the joint distribution of our training data and sample features and labels using that distribution and then train our model using that sampled data. At the end of the experiment we will try to draw some conclusion using the comparison of results.

```
In [1]: ## Importing relevant libraries

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
import MyLazyPredict
from sklearn.mixture import GaussianMixture
```

```
In [2]: ## Loading the dataset

housing = fetch_california_housing()
features = pd.DataFrame(housing.data, columns=housing.feature_names)
target = pd.DataFrame(housing.target, columns=housing.target_names)
housing = pd.concat([features, target], axis=1)
print(housing.shape)
housing.head()

(20640, 9)
```

```
Out[2]:
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	MedHouseVal
0	8.33	41.00	6.98	1.02	322.00	2.56	37.88	-122.23	4.53
1	8.30	21.00	6.24	0.97	2401.00	2.11	37.86	-122.22	3.58
2	7.26	52.00	8.29	1.07	496.00	2.80	37.85	-122.24	3.52
3	5.64	52.00	5.82	1.07	558.00	2.55	37.85	-122.25	3.41
4	3.85	52.00	6.28	1.08	565.00	2.18	37.85	-122.25	3.42

```
In [3]: ## Declaring features and target

X = housing.drop("MedHouseVal", axis=1)
y = housing["MedHouseVal"]
```

```
In [4]: ## Splitting the data into train and test datasets

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.8, random_state =42)
train_data = pd.concat([X_train, y_train], axis=1)
```

In [5]: *## visualizing train data distribution*

```
fig, axes = plt.subplots(3, 3, figsize=(8,8))
fig.suptitle('Distribution Plots of All Features', size=10)

sns.histplot(ax=axes[0,0], x=train_data["MedInc"], kde=True)
axes[0,0].set_xlabel('MedInc', size=8)

sns.histplot(ax=axes[0,1], x=train_data["HouseAge"], kde=True)
axes[0,1].set_xlabel('HouseAge', size=8)

sns.histplot(ax=axes[0,2], x=train_data["AveRooms"], kde=True)
axes[0,2].set_xlabel('AveRooms', size=8)

sns.histplot(ax=axes[1,0], x=train_data["AveBedrms"], kde=True)
axes[1,0].set_xlabel('AveBedrms', size=8)

sns.histplot(ax=axes[1,1], x=train_data["Population"], kde=True)
axes[1,1].set_xlabel('Population', size=8)

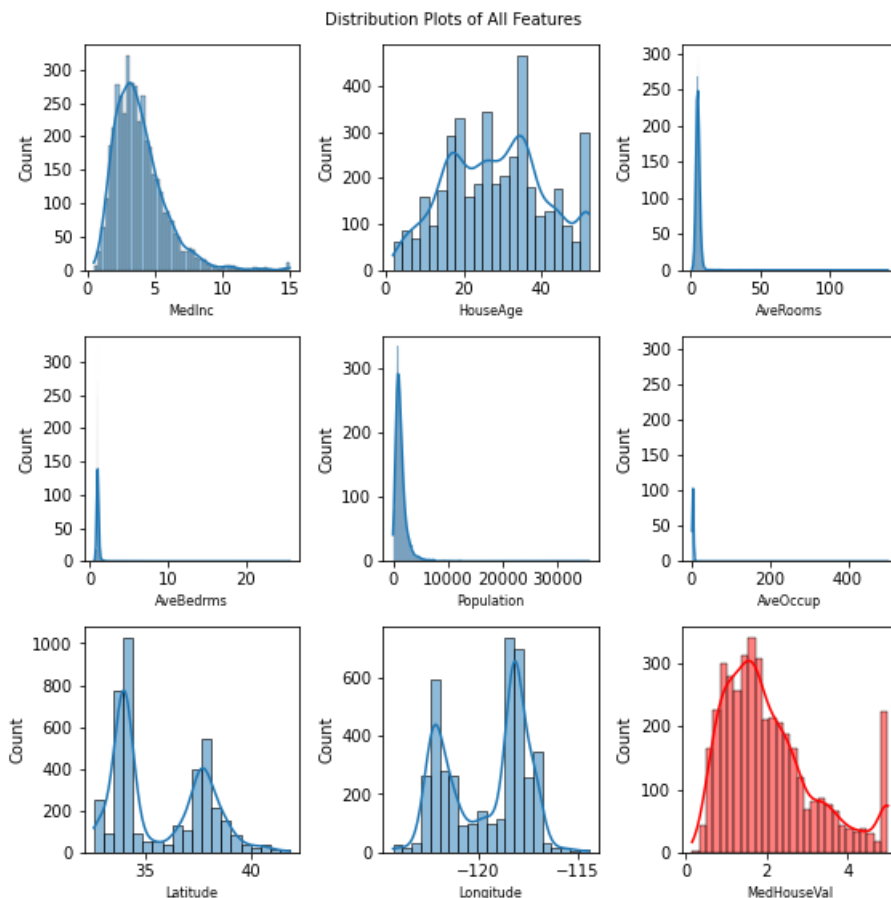
sns.histplot(ax=axes[1,2], x=train_data["AveOccup"], kde=True)
axes[1,2].set_xlabel('AveOccup', size=8)

sns.histplot(ax=axes[2,0], x=train_data["Latitude"], kde=True)
axes[2,0].set_xlabel('Latitude', size=8)

sns.histplot(ax=axes[2,1], x=train_data["Longitude"], kde=True)
axes[2,1].set_xlabel('Longitude', size=8)

sns.histplot(ax=axes[2,2], x=train_data["MedHouseVal"], kde=True, color='red')
axes[2,2].set_xlabel('MedHouseVal', size=8)

plt.tight_layout()
plt.show()
```



Here, we can see that there are normal distribution in the data. Also, some of the features have multimodal distributions.



In [8]: *## visualizing feature variables*

```
fig, axes = plt.subplots(3, 3, figsize=(8,8))
fig.suptitle('Distribution Plots of All Features', size=10)

sns.histplot(ax=axes[0,0], x=sampled_train_data["MedInc"], kde=True)
axes[0,0].set_xlabel('MedInc', size=8)

sns.histplot(ax=axes[0,1], x=sampled_train_data["HouseAge"], kde=True)
axes[0,1].set_xlabel('HouseAge', size=8)

sns.histplot(ax=axes[0,2], x=sampled_train_data["AveRooms"], kde=True)
axes[0,2].set_xlabel('AveRooms', size=8)

sns.histplot(ax=axes[1,0], x=sampled_train_data["AveBedrms"], kde=True)
axes[1,0].set_xlabel('AveBedrms', size=8)

sns.histplot(ax=axes[1,1], x=sampled_train_data["Population"], kde=True)
axes[1,1].set_xlabel('Population', size=8)

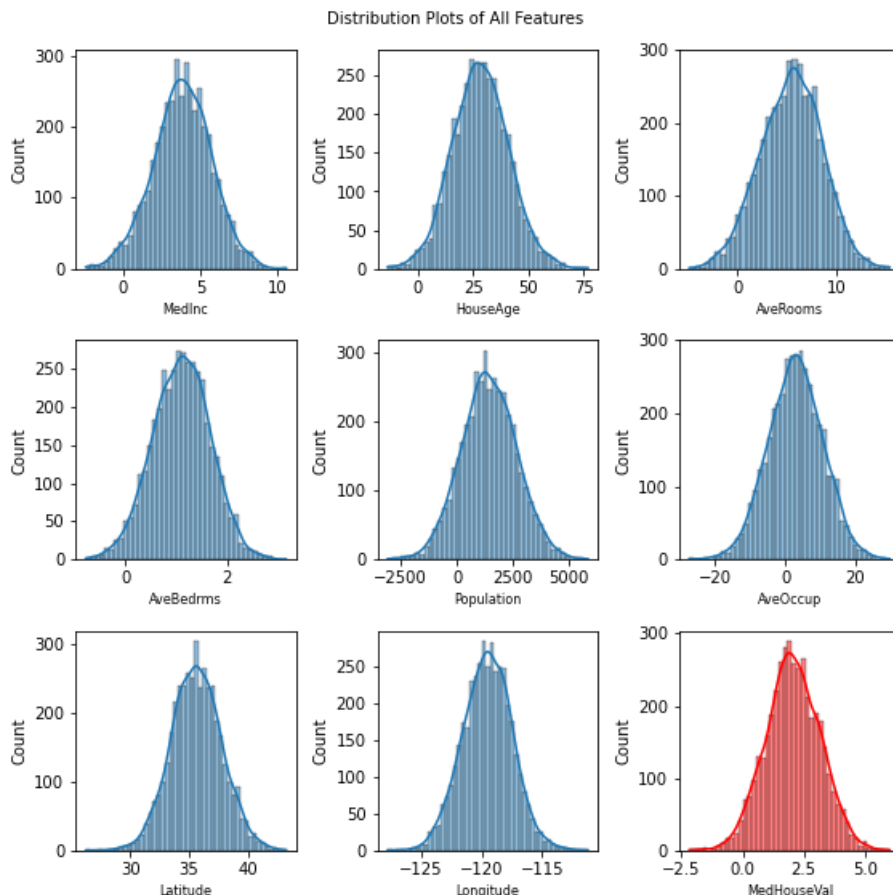
sns.histplot(ax=axes[1,2], x=sampled_train_data["AveOccup"], kde=True)
axes[1,2].set_xlabel('AveOccup', size=8)

sns.histplot(ax=axes[2,0], x=sampled_train_data["Latitude"], kde=True)
axes[2,0].set_xlabel('Latitude', size=8)

sns.histplot(ax=axes[2,1], x=sampled_train_data["Longitude"], kde=True)
axes[2,1].set_xlabel('Longitude', size=8)

sns.histplot(ax=axes[2,2], x=sampled_train_data["MedHouseVal"], kde=True, color='red')
axes[2,2].set_xlabel('MedHouseVal', size=8)

plt.tight_layout()
plt.show()
```



In [9]: *## Declaring Features and targets*

```
X_train_sampled = sampled_train_data.drop("MedHouseVal", axis=1)
y_train_sampled = sampled_train_data["MedHouseVal"]
```

Out[10]:

	Adjusted R-Squared	R-Squared	RMSE	Time Taken
Model				
LassoCV	0.60	0.60	0.72	0.09
LinearRegression	0.60	0.60	0.72	0.03
RidgeCV	0.60	0.60	0.72	0.03
Ridge	0.60	0.60	0.72	0.03
BayesianRidge	0.60	0.60	0.72	0.02
ElasticNetCV	0.60	0.60	0.72	0.15
LinearSVR	0.60	0.60	0.72	0.05
SGDRegressor	0.60	0.60	0.73	0.05
SVR	0.59	0.59	0.74	6.72
MLPRegressor	0.58	0.58	0.75	3.43
GradientBoostingRegressor	0.56	0.56	0.76	0.98
LGBMRegressor	0.56	0.56	0.77	0.19
XGBRegressor	0.54	0.54	0.78	0.58
RandomForestRegressor	0.54	0.54	0.79	2.44
BaggingRegressor	0.50	0.50	0.82	0.27
AdaBoostRegressor	0.47	0.47	0.84	0.42
KNeighborsRegressor	0.41	0.41	0.89	0.71
ElasticNet	0.17	0.17	1.05	0.02
DecisionTreeRegressor	0.12	0.12	1.08	0.06
Lasso	-0.00	-0.00	1.15	0.00

This time we have used all features and target variable to build our joint gaussian distribution. We will now build the mixture using only features.

The assumption, here, is that now the models will better predict as compared to the last experiment.

```
In [11]: gm = GaussianMixture(n_components=1, random_state=0).fit(housing)
sampled_train_data = pd.DataFrame(gm.sample(len(X_train))[0], columns=housing.columns)
sampled_train_data.head()
```

Out[11]:

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude	MedHouseVal
0	8.04	31.82	10.43	1.57	-572.20	-8.85	40.04	-123.53	3.31
1	3.56	30.21	6.64	1.38	960.51	-12.36	36.86	-121.28	1.39
2	6.74	37.11	6.77	0.94	1070.96	30.22	37.48	-120.06	1.95
3	5.13	11.35	6.20	1.17	1637.41	-14.01	35.51	-120.00	4.12
4	4.29	18.21	2.58	0.71	32.20	4.49	34.90	-120.76	3.23

In [12]: *## visualizing feature variables*

```
fig, axes = plt.subplots(3, 3, figsize=(8,8))
fig.suptitle('Distribution Plots of All Features', size=10)

sns.histplot(ax=axes[0,0], x=sampled_train_data["MedInc"], kde=True)
axes[0,0].set_xlabel('MedInc', size=8)

sns.histplot(ax=axes[0,1], x=sampled_train_data["HouseAge"], kde=True)
axes[0,1].set_xlabel('HouseAge', size=8)

sns.histplot(ax=axes[0,2], x=sampled_train_data["AveRooms"], kde=True)
axes[0,2].set_xlabel('AveRooms', size=8)

sns.histplot(ax=axes[1,0], x=sampled_train_data["AveBedrms"], kde=True)
axes[1,0].set_xlabel('AveBedrms', size=8)

sns.histplot(ax=axes[1,1], x=sampled_train_data["Population"], kde=True)
axes[1,1].set_xlabel('Population', size=8)

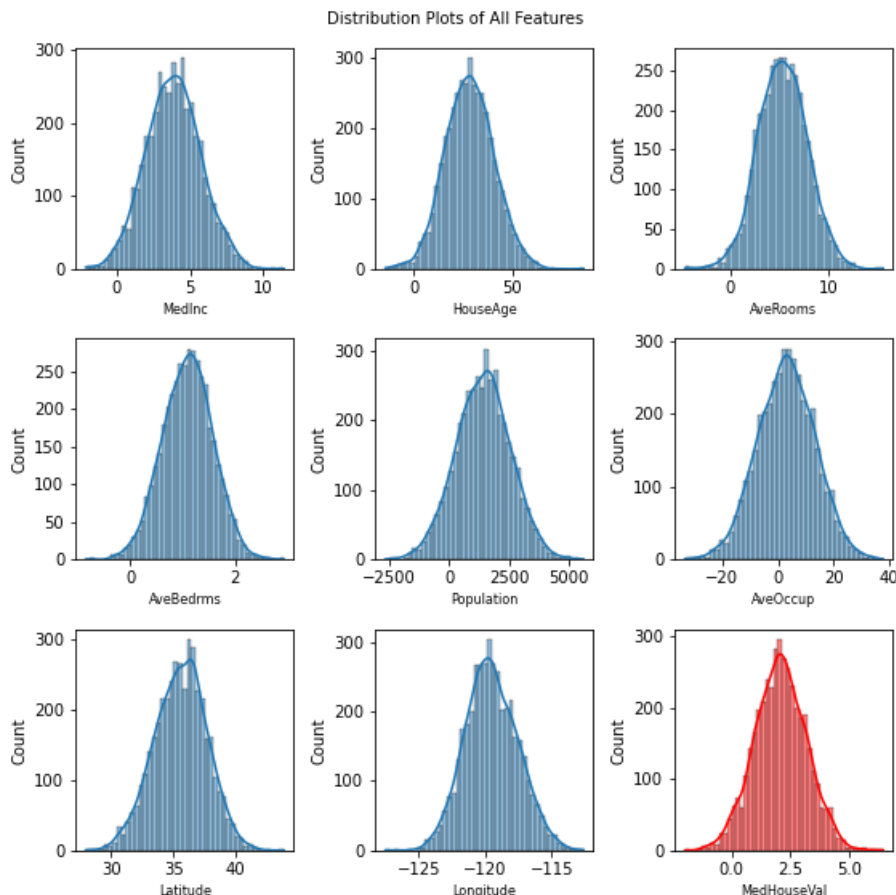
sns.histplot(ax=axes[1,2], x=sampled_train_data["AveOccup"], kde=True)
axes[1,2].set_xlabel('AveOccup', size=8)

sns.histplot(ax=axes[2,0], x=sampled_train_data["Latitude"], kde=True)
axes[2,0].set_xlabel('Latitude', size=8)

sns.histplot(ax=axes[2,1], x=sampled_train_data["Longitude"], kde=True)
axes[2,1].set_xlabel('Longitude', size=8)

sns.histplot(ax=axes[2,2], x=sampled_train_data["MedHouseVal"], kde=True, color='red')
axes[2,2].set_xlabel('MedHouseVal', size=8)

plt.tight_layout()
plt.show()
```



In [13]: *## Declaring Features and targets*

```
X_train_sampled = sampled_train_data.drop("MedHouseVal", axis=1)
y_train_sampled = sampled_train_data["MedHouseVal"]
```

Out[14]:

	Adjusted R-Squared	R-Squared	RMSE	Time Taken
Model				
SGDRegressor	0.61	0.61	0.72	0.02
LinearSVR	0.61	0.61	0.72	0.03
LinearRegression	0.61	0.61	0.72	0.03
Ridge	0.61	0.61	0.72	0.02
RidgeCV	0.61	0.61	0.72	0.02
BayesianRidge	0.61	0.61	0.72	0.02
LassoCV	0.61	0.61	0.72	0.09
ElasticNetCV	0.61	0.61	0.72	0.09
MLPRegressor	0.60	0.60	0.72	3.45
SVR	0.60	0.60	0.73	6.41
LGBMRegressor	0.57	0.57	0.75	0.19
GradientBoostingRegressor	0.56	0.56	0.77	0.94
RandomForestRegressor	0.54	0.54	0.78	2.43
XGBRegressor	0.53	0.53	0.79	0.59
BaggingRegressor	0.50	0.50	0.81	0.30
AdaBoostRegressor	0.46	0.46	0.84	0.33
KNeighborsRegressor	0.43	0.43	0.87	0.66
ElasticNet	0.18	0.18	1.05	0.02
DecisionTreeRegressor	0.11	0.11	1.09	0.05
Lasso	-0.00	-0.00	1.15	0.02

## Joint Distribution of Independent Features of the Data and Model Training

Out[15]:

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	6.36	18.07	9.91	1.68	3593.37	-3.71	40.29	-124.63
1	4.34	24.21	9.41	1.70	1313.79	1.54	38.18	-122.00
2	2.84	25.82	4.77	1.30	3263.80	1.70	29.43	-114.25
3	0.27	37.85	2.91	0.91	4210.79	5.14	37.98	-121.85
4	1.00	55.38	4.58	1.23	356.06	6.46	38.45	-122.01



In [16]: *## visualizing feature variables*

```
fig, axes = plt.subplots(3, 3, figsize=(8,8))
fig.suptitle('Distribution Plots of All Features', size=10)

sns.histplot(ax=axes[0,0], x=sampled_train_data["MedInc"], kde=True)
axes[0,0].set_xlabel('MedInc', size=8)

sns.histplot(ax=axes[0,1], x=sampled_train_data["HouseAge"], kde=True)
axes[0,1].set_xlabel('HouseAge', size=8)

sns.histplot(ax=axes[0,2], x=sampled_train_data["AveRooms"], kde=True)
axes[0,2].set_xlabel('AveRooms', size=8)

sns.histplot(ax=axes[1,0], x=sampled_train_data["AveBedrms"], kde=True)
axes[1,0].set_xlabel('AveBedrms', size=8)

sns.histplot(ax=axes[1,1], x=sampled_train_data["Population"], kde=True)
axes[1,1].set_xlabel('Population', size=8)

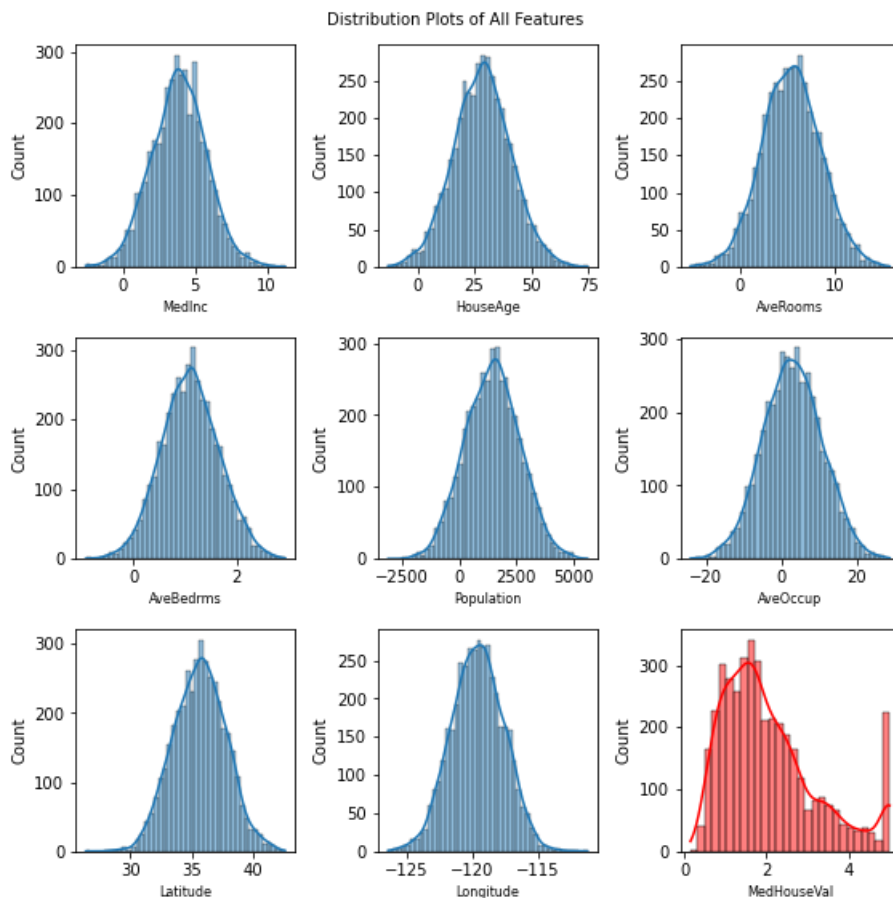
sns.histplot(ax=axes[1,2], x=sampled_train_data["AveOccup"], kde=True)
axes[1,2].set_xlabel('AveOccup', size=8)

sns.histplot(ax=axes[2,0], x=sampled_train_data["Latitude"], kde=True)
axes[2,0].set_xlabel('Latitude', size=8)

sns.histplot(ax=axes[2,1], x=sampled_train_data["Longitude"], kde=True)
axes[2,1].set_xlabel('Longitude', size=8)

sns.histplot(ax=axes[2,2], x=train_data["MedHouseVal"], kde=True, color='red')
axes[2,2].set_xlabel('MedHouseVal', size=8)

plt.tight_layout()
plt.show()
```



```
In [17]: reg = MyLazyPredict.LazyRegressor(verbose=0, ignore_warnings=False, custom_metric=None)
models, predictions = reg.fit(sampled_train_data, X_test, y_train, y_test)
models
```

[illegible]

Out[17]:

	Adjusted R-Squared	R-Squared	RMSE	Time Taken
Model				
ElasticNet	-0.00	-0.00	1.15	0.02
ElasticNetCV	-0.00	-0.00	1.15	0.09
Lasso	-0.00	-0.00	1.15	0.02
LassoCV	-0.00	-0.00	1.15	0.09
BayesianRidge	-0.00	-0.00	1.15	0.02
PoissonRegressor	-0.01	-0.01	1.16	0.03
GradientBoostingRegressor	-0.02	-0.02	1.16	1.02
SGDRegressor	-0.03	-0.03	1.17	0.02
RidgeCV	-0.03	-0.03	1.17	0.02
Ridge	-0.03	-0.03	1.17	0.02
LinearRegression	-0.03	-0.03	1.17	0.02
MLPRegressor	-0.04	-0.04	1.18	3.30
AdaBoostRegressor	-0.05	-0.05	1.18	0.09
RandomForestRegressor	-0.06	-0.06	1.18	3.37
LinearSVR	-0.06	-0.06	1.18	0.03
SVR	-0.13	-0.13	1.23	6.74
LGBMRegressor	-0.14	-0.14	1.23	0.19
XGBRegressor	-0.15	-0.15	1.24	0.56
BaggingRegressor	-0.15	-0.15	1.24	0.36
KNeighborsRegressor	-0.17	-0.17	1.25	0.65
DecisionTreeRegressor	-0.84	-0.84	1.56	0.09

Here, we can see that no model has able to capture the pattern (because, there are not any) as the joint distribution's objective has been compromised when we left the target variable out.

### Conclusion:

In above four experiments, we have seen that the best way of learning the joint distribution is to learn it with all the features and target variables. Also, we have seen, how incorporating the test data can enhance the models' performances. The last experiment was a total failure in terms of the performances as it on the backend did not incorporate the target variable while learning the features.

Learning joint gaussian distribution might come handy in different approaches. It can help us track the model performance and see where is it drifting if performing bad for newer records and why.

Also, many times we do not have enough data to learn complex machine learning algorithms, there it can help us generated through sampling enough data to train our models.