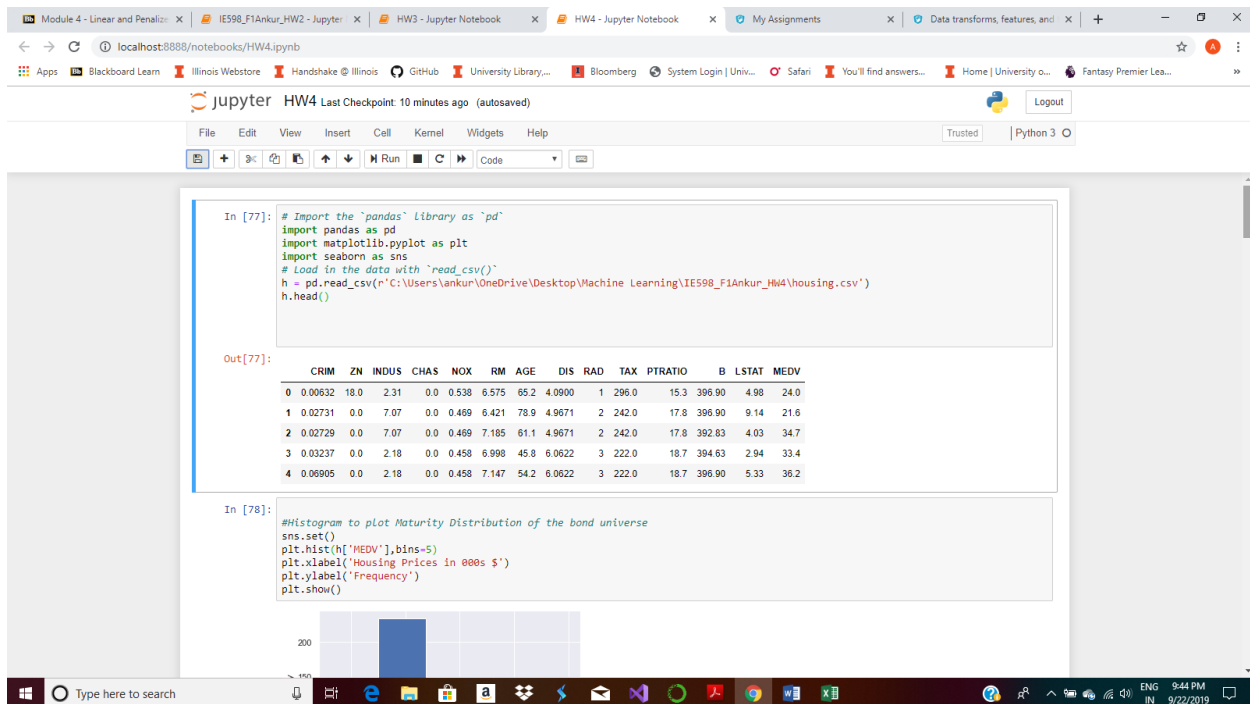


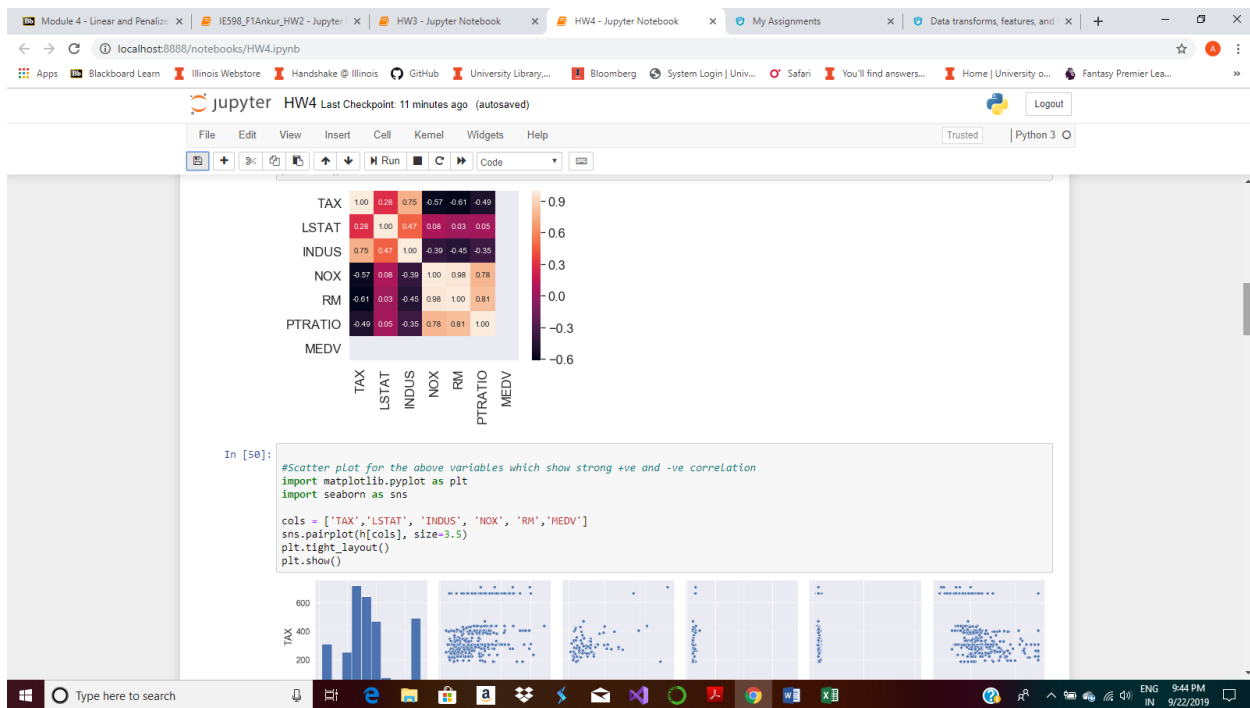
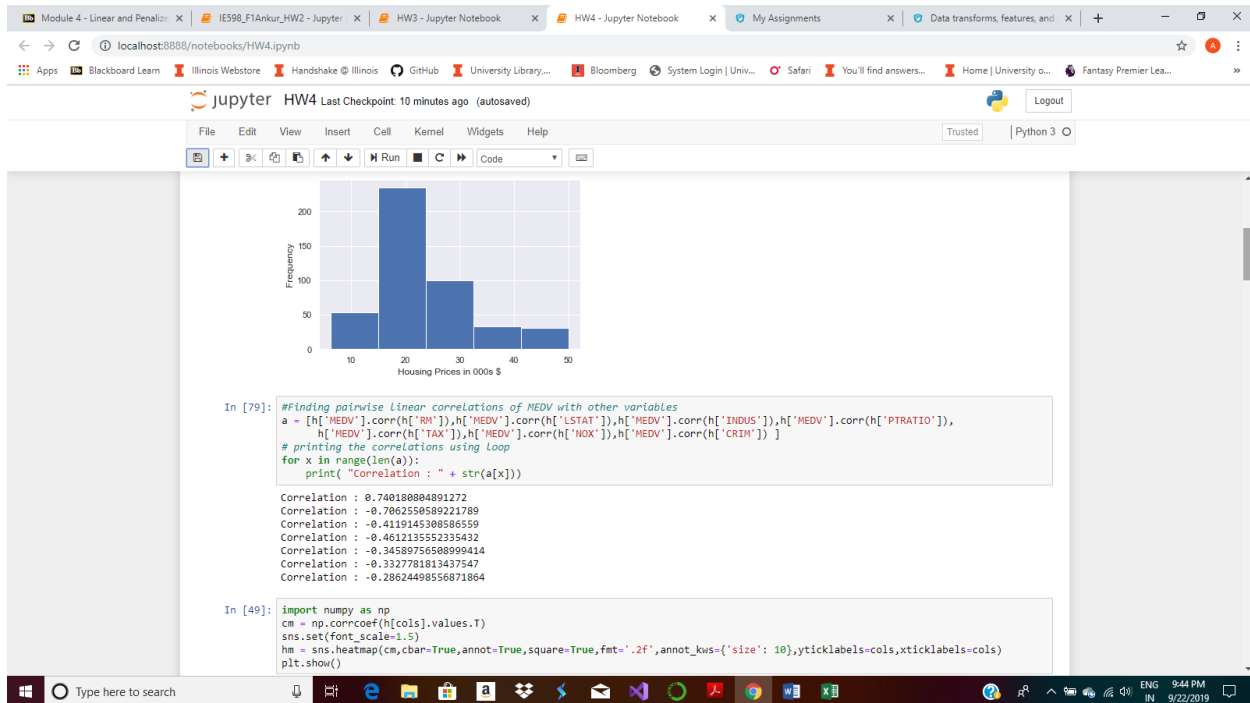
Ankur Mukherjee (ankurm3)

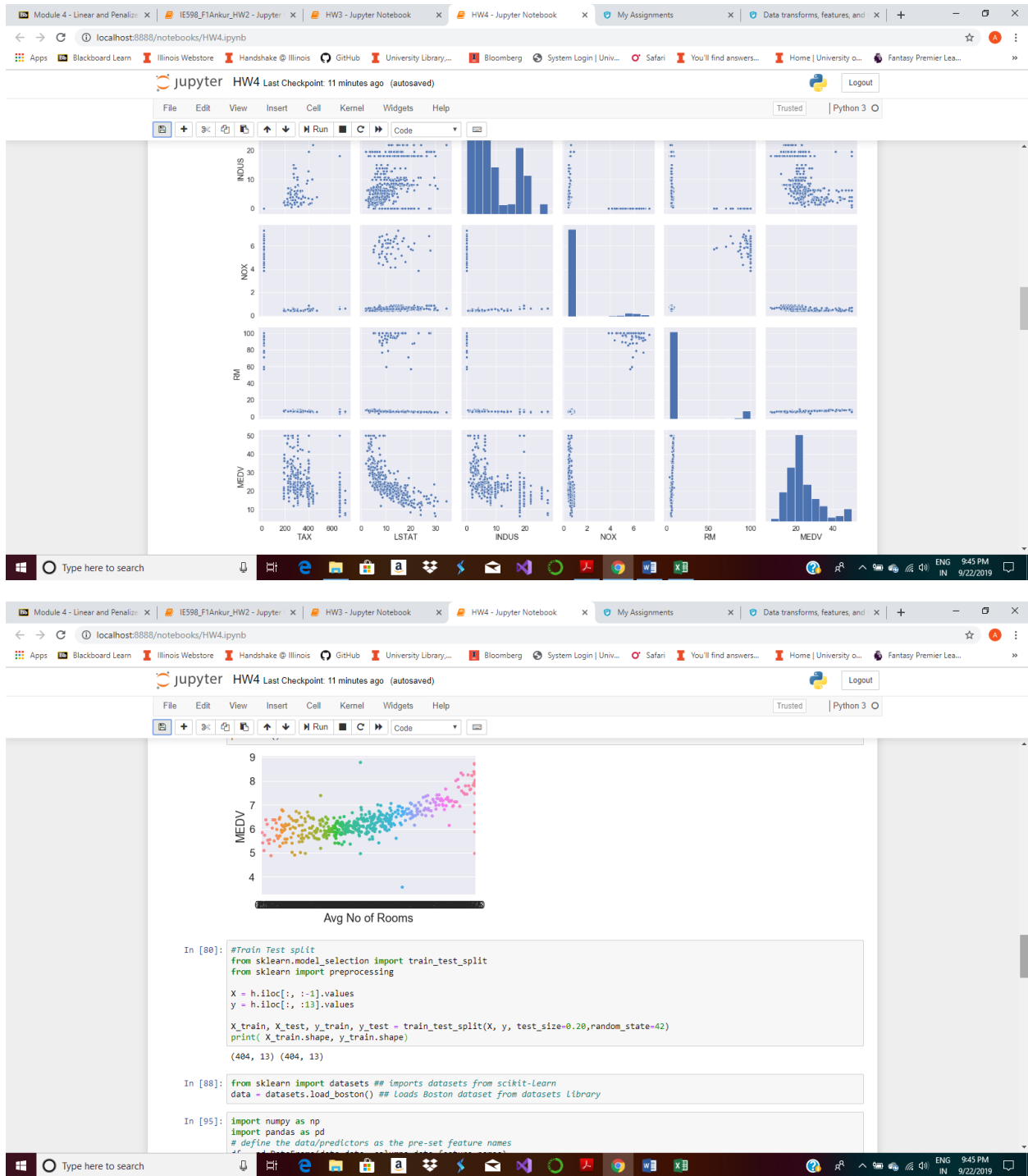
IE598 MLF F18

Module 4 Homework (Regression)

Part 1: Exploratory Data Analysis







Part 2: Linear regression

Jupyter HW4 Last Checkpoint: 21 minutes ago (autosaved)

```

In [118]: import statsmodels.api as sm

X = df[["RM", "LSTAT", "INDUS", "NOX", "TAX"]]
y = target["MEDV"]

# Note the difference in argument order
model = sm.OLS(y, X).fit()
predictions = model.predict(X) # make the predictions by the model

# Print out the statistics
model.summary()

```

Out[118]:

OLS Regression Results

Dep. Variable:	MEDV	R-squared (uncentered):	0.950			
Model:	OLS	Adj. R-squared (uncentered):	0.950			
Method:	Least Squares	F-statistic:	1907.			
Date:	Sun, 22 Sep 2019	Prob (F-statistic):	2.47e-323			
Time:	20:55:19	Log-Likelihood:	-1574.8			
No. Observations:	506	AIC:	3160.			
DF Residuals:	501	BIC:	3181.			
DF Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
RM	4.9800	0.194	25.672	0.000	4.599	5.361
LSTAT	-0.5914	0.047	-12.570	0.000	-0.684	-0.499
INDUS	0.0241	0.063	0.384	0.701	-0.099	0.148
NOX	3.2593	3.352	0.972	0.331	-3.327	9.845
TAX	-0.0082	0.002	-3.759	0.000	-0.012	-0.004

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 6.11e+03. This might indicate that there are strong multicollinearity or other numerical problems.

```

In [117]: #Mean Squared Error
import numpy as np
from sklearn.metrics import mean_squared_error
np.sqrt(mean_squared_error(y, df[["RM"]]))

```

Out[117]: 18.43763831298453

Part 3.1: Ridge regression

```
Out[117]: 18.43763831290453

In [134]: #Ridge regression
from sklearn.linear_model import Ridge
from sklearn.datasets import load_boston
import numpy as np
# Load data
boston = load_boston()
X = boston.data
y = boston.target

# Create ridge regression with an alpha value
regr = Ridge(alpha=0.5)

# Fit the linear regression
regr.fit(X_std, y)
Ridge(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=None,
normalize=False, random_state=None, solver='auto', tol=0.001)
#R^2
regr.score(X_std,y)

Out[134]: 0.7406395408018083

In [148]: #R^2 is maximized at alpha = 0.1
regr = Ridge(alpha=0.1)

# Fit the linear regression
regr.fit(X_std, y)
#R^2
regr.score(X_std,y)

Out[148]: 0.7406426322193994

In [144]: #Lasso
from sklearn.linear_model import Lasso
from sklearn.datasets import load_boston

# Load data
boston = load_boston()
X = boston.data
y = boston.target

# Create ridge regression with an alpha value
regr = Lasso(alpha=0.5)

# Fit the linear regression
regr.fit(X_std, y)

#R^2
regr.score(X_std,y)

Out[144]: 0.6913555608028645

In [149]: # Create ridge regression with max alpha value
regr = Lasso(alpha=0.1)

# Fit the linear regression
regr.fit(X_std, y)

#R^2
regr.score(X_std,y)

Out[149]: 0.7353093744353656

In [ ]: print("My name is Ankur Mukherjee")
print("My NetID is: ankurm3")
print("I hereby certify that I have read the University policy on Academic Integrity and that I am not in violation.")
```

Coefficient of Determination(R^2) seems to max out as alpha decreases to 0.1

Part 3.2: LASSO regression

```
In [144]: #Lasso
from sklearn.linear_model import Lasso
from sklearn.datasets import load_boston
import numpy as np
# Load data
boston = load_boston()
X = boston.data
y = boston.target

# Create ridge regression with an alpha value
regr = Lasso(alpha=0.5)

# Fit the linear regression
regr.fit(X_std, y)

#R^2
regr.score(X_std,y)

Out[144]: 0.6913555608028645

In [149]: # Create ridge regression with max alpha value
regr = Lasso(alpha=0.1)

# Fit the linear regression
regr.fit(X_std, y)

#R^2
regr.score(X_std,y)

Out[149]: 0.7353093744353656

In [ ]: print("My name is Ankur Mukherjee")
print("My NetID is: ankurm3")
print("I hereby certify that I have read the University policy on Academic Integrity and that I am not in violation.")
```

Coefficient of Determination(R^2) seems to max out as alpha decreases to 0.1

Part 4: Conclusions

- 1) The linear regression model has a $\text{adj } R^2$ of 0.95. This means that 95% of the variation in MEDV is explained by the combination of RM, INDUS, NOX, LSTAT and TAX variables
- 2) Some of the coefficients have low t-stats but The f-statistics is high which means together they can explain the variation but not individually. This indicates **multicollinearity**
- 3) MEDV is not normally distributed with high kurtosis(fat tails), right positive skewed with a few high outliers – few houses costing more.

Part 5: Appendix

https://github.com/ankurmukherjeeuiuc/IE598_F1Ankur_HW4