

▼ Homework 2

▼ Description:

Columns represent (in order): 1) Company where they work 2) Job title 3) Office location 4) Total annual compensation (in \$) 5) *Base salary* (in \$) 6) Value of stock grants (in \$) 7) *Bonus payments* (in \$) 8) Years of relevant experience (in years) 9) Time with this company (in years) 10) Gender (self-reported) 11) Terminal Degree is Masters (1 = yes) 12) Terminal Degree is Bachelors (1 = yes) 13) Terminal Degree is Doctorate (1 = yes) 14) Terminal Degree is High School (1 = yes) 15) Terminal Degree is some college (1 = yes) 16) Self-identifies as Asian (1 = yes) 17) Self-identifies as White (1 = yes) 18) Self-identifies as Multi-Racial (1 = yes) 19) Self-identifies as Black (1 = yes) 20) Self-identifies as Hispanic (1 = yes) 21) Race as a qualitative variable 22) Education as a qualitative variable 23) Age (in years) 24) Height (in inches) 25) Zodiac sign (Tropical calendar, 1 = Aries, 12 = Pisces, with everything else in between) 26) SAT score 27) GPA We/you will want to use most of these variables in prediction models. This data is self-reported, so sometimes it will be missing if the person (for whatever reason) did not provide this information. For instance, the information on education (variables 11-15) is only meaningfully interpretable for any given row, if the corresponding value in variable 22 is not "NA". NA indicates missing data. The same is true for variables 16 to 21.

▼ Loading the data:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import MinMaxScaler
4 from sklearn.preprocessing import LabelEncoder
5 import matplotlib.pyplot as plt
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import r2_score
8 from sklearn.linear_model import Ridge
9 from sklearn.model_selection import train_test_split
10 from sklearn.metrics import mean_squared_error
11
12 import warnings
13 warnings.filterwarnings('ignore')
14 data = pd.read_csv("techSalaries2017.csv")
```

▼ Data Cleaning/Engineering:

```
1 data.head(10)
2 data = data.drop(['basesalary', 'stockgrantvalue', 'bonus', 'Zodiac'], axis = 1)
3 main_data = data
```

```
1 print(data.gender.value_counts())
2 print(data.Race.value_counts())
3 data['company'] = data['company'].str.upper()
4 data['title'] = data['title'].str.upper()
5 data['location'] = data['location'].str.upper()
6
7 # dropping rows with nan for company name,
8 # I believe this predictor influences the predictions a lot
9 data = data[data['company'].notna()]
10
11
12
```

```
Male          35702
Female        6999
Other          400
Title: Senior Software Engineer    1
Name: gender, dtype: int64
Asian        11772
White         8032
Hispanic     1129
Two Or More   804
Black         690
Name: Race, dtype: int64
```

```
1 data.sample(5)
```

```
2
```

	company	title	location	totalyearlycompensation	years
12066	VANGUARD	SOFTWARE ENGINEER	MALVERN, PA	80000	
27940	ROBINHOOD	SOFTWARE ENGINEER	SAN FRANCISCO, CA	194000	
22031	AMAZON	SOFTWARE ENGINEER	SAN FRANCISCO, CA	330000	
17973	BLOOMBERG	SOFTWARE ENGINEERING MANAGER	NEW YORK, NY	357000	
37544	ZILLOW	SOFTWARE ENGINEERING MANAGER	IRVINE, CA	260000	

5 rows x 23 columns

▼ Encoding all categorical variables:

```
1 le = LabelEncoder()
2 data['company'] = le.fit_transform(data['company'])
3 data['title'] = le.fit_transform(data['title'])
4 data['location'] = le.fit_transform(data['location'])
5
```

▼ Rescaling features

```

1 # normalizing the outcome variable
2 # Assuming the column to be normalized is stored in a NumPy array X
3 column_to_normalize = np.array(data['totalyearlycompensation']) # Select the
4 scaler = MinMaxScaler() # Create a MinMaxScaler object
5 outcome = scaler.fit_transform(column_to_normalize.reshape(-1, 1))
6 data['totalyearlycompensation']=outcome
7
8 '''
9 I had previously used these 3 variables until I read the hints,
10 So I will abstain from them for the rest of the assignment.
11 '''
12 # column_to_normalize = np.array(data['company']) # Select the column you wa
13 # data['company'] = scaler.fit_transform(column_to_normalize.reshape(-1, 1))
14
15 # column_to_normalize = np.array(data['title']) # Select the column you want
16 # data['title'] = scaler.fit_transform(column_to_normalize.reshape(-1, 1))
17
18 # column_to_normalize = np.array(data['location']) # Select the column you w
19 # data['location'] = scaler.fit_transform(column_to_normalize.reshape(-1, 1))
20
21 column_to_normalize = np.array(data['yearsofexperience']) # Select the colum
22 data['yearsofexperience'] = scaler.fit_transform(column_to_normalize.reshape(
23
24 column_to_normalize = np.array(data['yearsatcompany']) # Select the column y
25 data['yearsatcompany'] = scaler.fit_transform(column_to_normalize.reshape(-1,
26 # column_to_normalize = np.array(data['gender']) # Select the column you wan
27 # data['gender'] = scaler.fit_transform(column_to_normalize.reshape(-1, 1))
28
29 predictors = data.drop(['totalyearlycompensation',"company",'Race_White','Rac
30

```

```
1 data.head(5)
```

	company	title	location	totalyearlycompensation	yearsofexperience	yearsatcompany
0	722	8	767	0.023541	0.021739	0.021739
1	331	11	823	0.018109	0.072464	0.072464
2	47	8	859	0.060362	0.115942	0.115942
3	71	12	909	0.072837	0.101449	0.101449
4	640	11	610	0.029577	0.072464	0.072464

5 rows x 23 columns

1. Using multiple linear regression: What is the best predictor of
- ▼ total annual compensation, how much variance is explained by this predictor vs. the full multiple regression model?

Answer:

On the preprocessing

- The variables 5 to 7 sounds to be directly related to the compensation and therefore might be irrelevant to predicting the total compensation. Example: someone with high total compensation might be more likely to be in a position with higher base salary, stock grants and bonuses.
- Columns like Height, age and zodiac were disregarded because either they make no sense on the relation with yearly compensation or because they might be correlated with other variable already being taken into account. Example: Years of experience may be very correlated with the age -the older you are more experience you have-. Later on you the correlations table you will find that it's true.
- Other columns were removed because they had already been turned into categorical variables such as education and race.

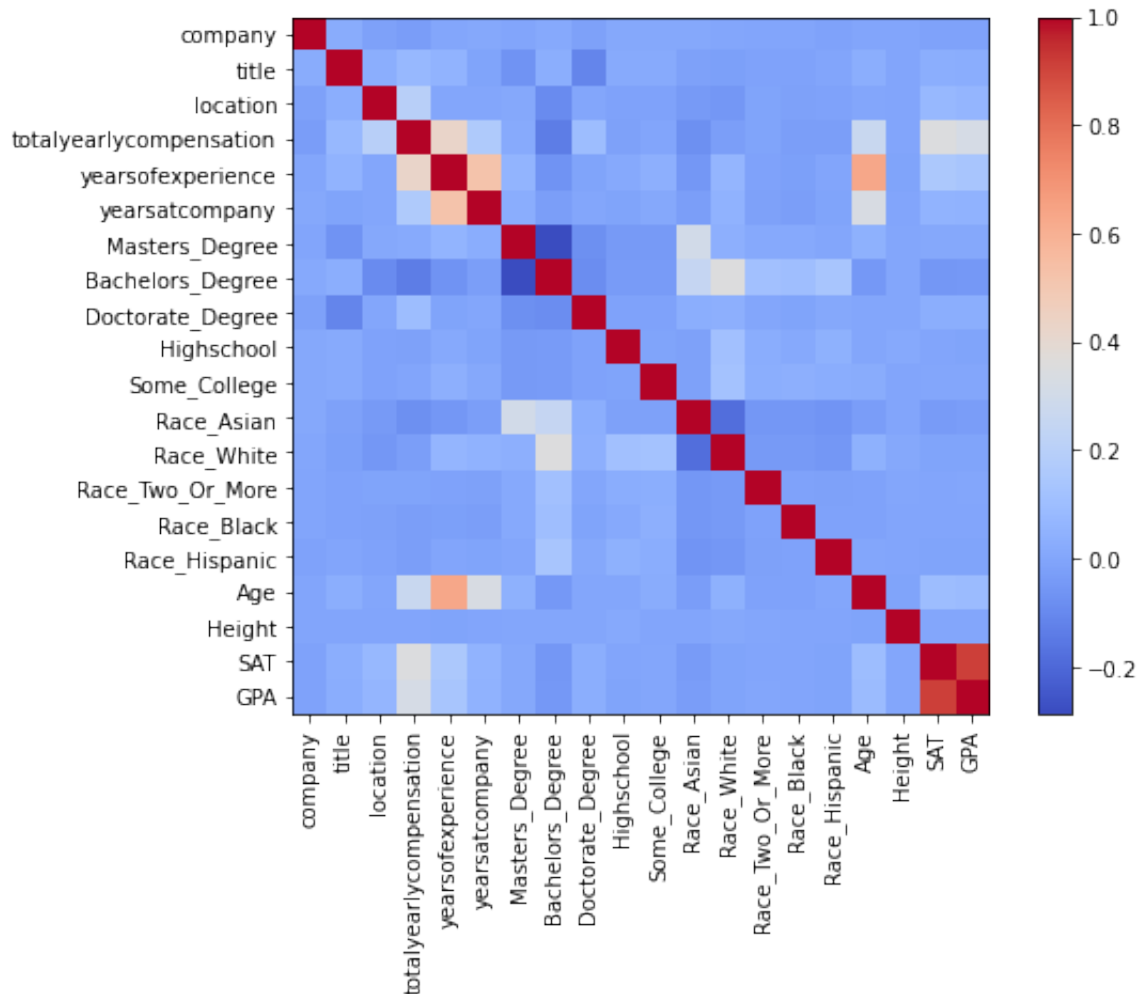
▼ Checking for correlations

- Columns that were found to be highly correlated were removed.
- The correlation table and heatmap were useful for picking a best predictor. The best predictor chosen by me was the `yearsofexperience` column because it has the largest correlation coefficient with our target column `totalyearlycompensation`.
- After looking for the least correlated variables for a second time, I also dropped columns associated with race. After that I saw a slightly increase on the R^2 of both models. This was not documented but I observed it through a process of trial and error.

```

1 df = pd.DataFrame(data)
2 plt.figure(figsize=(8, 6))# get the correlation matrix
3 corr = df.corr()
4
5 # plot the correlation matrix using matplotlib
6 plt.imshow(corr, cmap='coolwarm', interpolation='nearest')
7 plt.colorbar()
8 plt.xticks(range(len(corr)), corr.columns, rotation=90)
9 plt.yticks(range(len(corr)), corr.columns)
10 plt.show()

```



```
1 corr.totalyearlycompensation.sort_values()
```

```
Bachelors_Degree      -0.139833
Race_Asian             -0.080103
company               -0.030145
Race_White            -0.027791
Race_Black            -0.026748
Race_Hispanic         -0.026024
Highschool            -0.014836
Race_Two_Or_More      -0.009629
Some_College          -0.003384
Height                -0.001359
Masters_Degree         0.018288
title                 0.073268
Doctorate_Degree       0.097173
yearsatcompany         0.165774
location              0.196606
Age                   0.264170
GPA                   0.319253
SAT                   0.349962
yearsofexperience      0.422848
totalyearlycompensation 1.000000
Name: totalyearlycompensation, dtype: float64
```


▼ Regressions

- Now comparing the 2 regressions. By definition R-squared (R^2) is a value between 0 and 1 that represents the proportion of the variation in the dependent variable explained by the independent variables in a regression model. By comparing the 2 models R^2 we can observe that the simple regression model using our best predictor `yearsofexperience` has explained around 18.5% of the variance on the dependable variable `totalyearlycompensation`. While our multiple regression model with handpicked predictors, have achieved an R^2 of about .33 which means that this model has explained around 33% of the variance on our dependable variable.
- One could argue that our multiple regression model is over fitting by having multiple predictors, but our RMSE has shown that not to be true.

Simple Linear Regression:

R-squared: 0.1846469238598084 RMSE: 0.0006212221234210354

- Edited after review

Multiple Linear Regression:

Containing predictors title, company, and location R-squared: 0.3259203094713786 RMSE: 0.0005135851313488935

without predictors title, company, and location R-squared: 0.2952100301269125 RMSE: 0.0005369834669945201

- On the edited part:
 - It is plausible that even though the model with the 3 extra categorical predictors is better, it may also be true that it over fits to the training data.

```

1 xTrain, xTest, yTrain, yTest = train_test_split(np.array(data['yearsofexperie
2
3 # create a regression object
4 model = LinearRegression()
5
6 # fit the model to your data
7 model.fit(xTrain, yTrain)
8
9 # predict the dependent variable
10 y_pred = model.predict(xTest)
11
12 # calculate R-squared
13 r_squared = r2_score(yTest, y_pred)
14 rmse = mean_squared_error(yTest,y_pred)
15
16 print("R-squared:", r_squared,"RMSE:",rmse )

```

R-squared: 0.1846469238598084 RMSE: 0.0006212221234210354

```

1 # create a regression object
2 xTrain, xTest, yTrain, yTest = train_test_split(predictors, outcome.reshape(-
3
4 model = LinearRegression()
5
6 # fit the model to your data
7 model.fit(xTrain, yTrain)
8
9 # predict the dependent variable
10 y_pred = model.predict(xTest)
11
12 # calculate R-squared
13 r_squared = r2_score(yTest, y_pred)
14 rmse = mean_squared_error(yTest,y_pred)
15
16 print("R-squared:", r_squared,"RMSE:",rmse )

```

R-squared: 0.2952100301269125 RMSE: 0.0005369834669945201

2. Using ridge regression to do the same as in 1): How does the
- ▼ model change or improve compared to OLS? What is the optimal lambda?

Answer:

Report:

Simple Linear Regression:

RMSE = 0.0006643937370673719 $R^2 = 0.17520990973078954$

Optimal lambda: -5.199999999999999

Multiplre Linear Regression:

Optimal lambda: -4.074999999999999

RMSE = 0.0005535932718914939 $R^2 = 0.3127595592468504$

Both models did change, but for the worst. I continued using the same random seed for the process of train split, so the models would be trained and the lambdas would be test on the same data as our original models on the question 1.

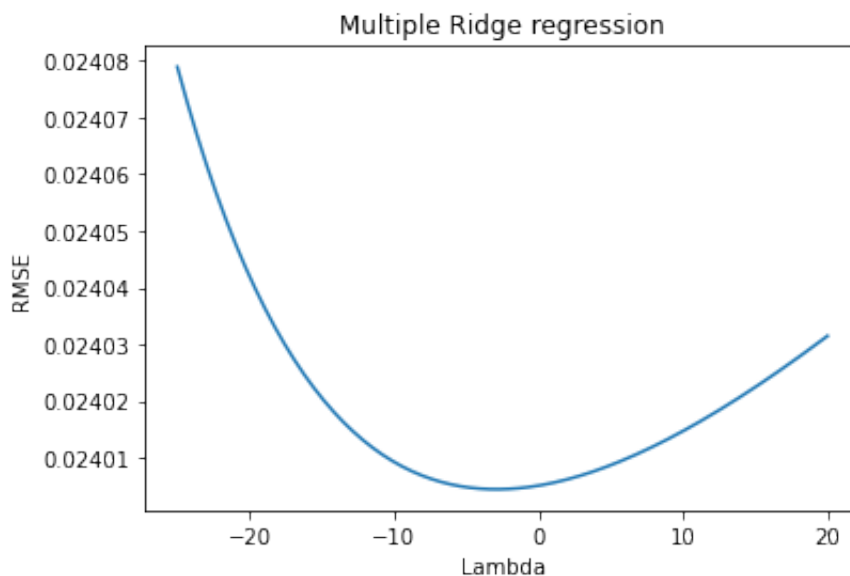
Still, both Ridge Regression models are worst by a margin of ~1%. Meaning that both regularized models explain about ~1% less then our un-regularized models.

▼ Multiple Ridge Regression

```

1
2 # Init parameters:
3 seed = 0
4 xTrain, xTest, yTrain, yTest = train_test_split(predictors, outcome.reshape(-
5 lambdas = np.linspace(-25,20,201)
6 cont = np.empty([len(lambdas),2])*np.NaN # [lambda error]
7 for ii in range(len(lambdas)):
8     ridgeModel = Ridge(alpha=lambdas[ii]).fit(xTrain, yTrain)
9     cont[ii,0] = lambdas[ii]
10    error = mean_squared_error(yTest,ridgeModel.predict(xTest),squared=False)
11    cont[ii,1] = error
12
13 plt.plot(cont[:,0],cont[:,1])
14 plt.xlabel('Lambda')
15 plt.ylabel('RMSE')
16 plt.title('Multiple Ridge regression')
17 plt.show()
18 optimal_lambda_multiple = lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))]
19 print('Optimal lambda:',optimal_lambda_multiple)
20 # multiple Ridge regression
21
22

```



Optimal lambda: -2.9499999999999993

```
1 # multiple Ridge regression
2 ridgeModel = Ridge(alpha=optimal_lambda_multiple).fit(xTrain, yTrain)
3 error = mean_squared_error(yTest,ridgeModel.predict(xTest))
4 r_squared = r2_score(yTest,ridgeModel.predict(xTest))
5 print(error,r_squared)
6
```

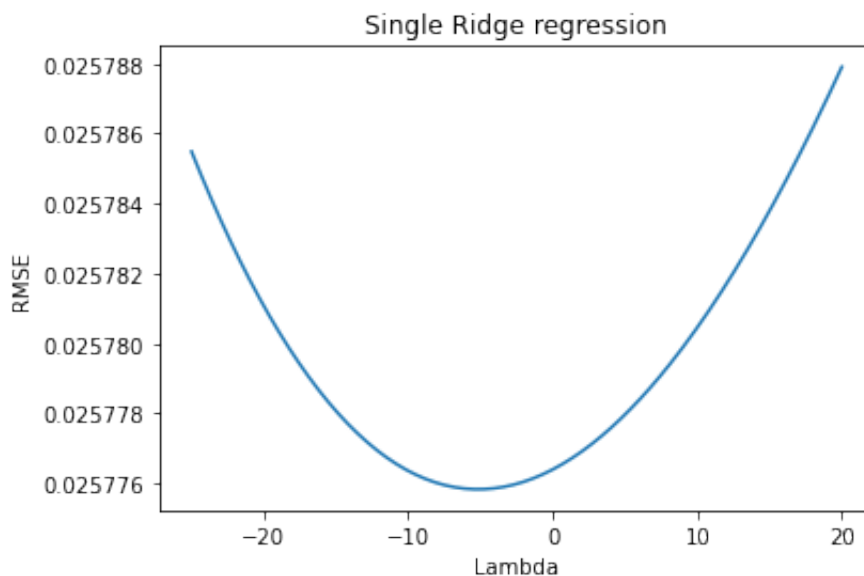
```
0.0005762144715295082 0.28467720348324754
```

▼ Single Ridge Regression

```

1 # Init parameters:
2 xTrain, xTest, yTrain, yTest = train_test_split(np.array(data['yearsofexperie
3 lambdas = np.linspace(-25,20,201)
4 cont = np.empty([len(lambdas),2])*np.NaN # [lambda error]
5 for ii in range(len(lambdas)):
6     ridgeModel = Ridge(alpha=lambdas[ii]).fit(xTrain, yTrain)
7     cont[ii,0] = lambdas[ii]
8     error = mean_squared_error(yTest,ridgeModel.predict(xTest),squared=False)
9     cont[ii,1] = error
10
11 plt.plot(cont[:,0],cont[:,1])
12 plt.xlabel('Lambda')
13 plt.ylabel('RMSE')
14 plt.title('Single Ridge regression')
15 plt.show()
16 optimal_lambda_single = lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))]
17
18 print('Optimal lambda:',optimal_lambda_single)

```



Optimal lambda: -5.199999999999999

```

1 # single Ridge regression
2 ridgeModel = Ridge(alpha=optimal_lambda_single).fit(xTrain, yTrain)
3 error = mean_squared_error(yTest,ridgeModel.predict(xTest),squared=True)
4 r_squared = r2_score(yTest,ridgeModel.predict(xTest))
5 print(error,r_squared)

```

0.0006643937370673719 0.17520990973078954

3. Using Lasso regression to do the same as in 1): How does the
- ▼ model change now? How many of the predictor betas are shrunk to exactly 0? What is the optimal lambda now?

Answer:

Both models now are actually performing worse. Probably due to a bad seed for the random state of the train_test_split function. on the multiple regression 9 of our predictor's betas were shrunk.

For multiple regression:

Optimal lambda: 0.5500000000000007

For Simple Linear regression:

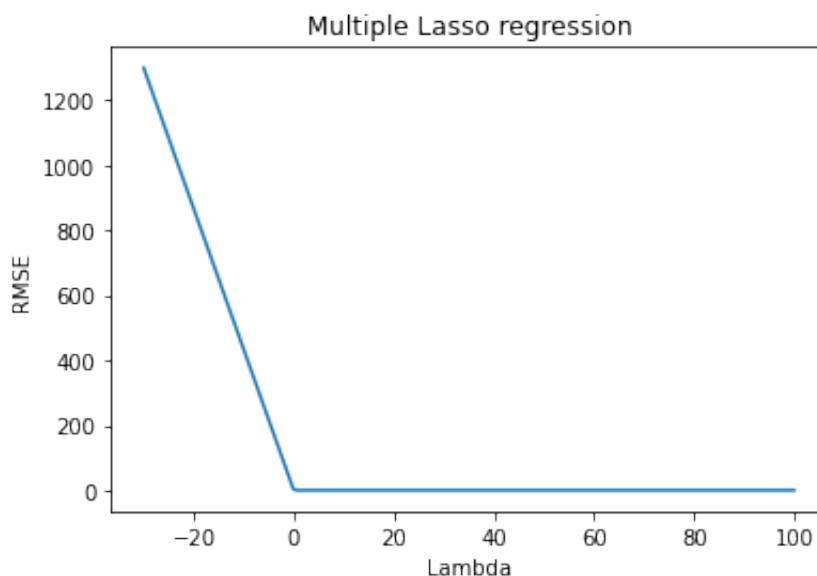
Optimal lambda: 0.13422818791946156

- ▼ Multiple Lasso Regression

```

1 ### 5. Now do the same thing--but with lasso regression
2 import warnings
3 warnings.filterwarnings('ignore') # Just to ignore warnings that might be thr
4
5 # Load libraries:
6 from sklearn.linear_model import Lasso
7
8 # Init parameters:
9 xTrain, xTest, yTrain, yTest = train_test_split(predictors, outcome.reshape(-
10 lambdas = np.linspace(-30,100,201)
11
12 cont = np.empty([len(lambdas),2])*np.NaN # [lambda error]
13
14 for ii in range(len(lambdas)):
15     lassoModel = Lasso(alpha=lambdas[ii]).fit(xTrain, yTrain)
16     cont[ii,0] = lambdas[ii]
17     error = mean_squared_error(yTest, lassoModel.predict(xTest), squared=False)
18     cont[ii,1] = error
19
20 plt.plot(cont[:,0], cont[:,1])
21 plt.xlabel('Lambda')
22 plt.ylabel('RMSE')
23 plt.title('Multiple Lasso regression')
24 plt.show()
25 optimal_lambda_multiple = lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))]
26 print('Optimal lambda:', lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))])
27

```



Optimal lambda: 0.5500000000000007


```

1
2 #multiple lasso
3 lassoModel = Lasso(alpha=optimal_lambda_multiple).fit(xTrain, yTrain)
4 error = mean_squared_error(yTest,lassoModel.predict(xTest),squared=False)
5 r_squared = r2_score( yTest,lassoModel.predict(xTest))
6 print(error,r_squared)
7 print(lassoModel.coef_)

0.4773686840218829 0.0884726320371435
[ 0.          0.          0.         -0.          0.         -0.
 -0.          0.0050059  -0.          0.00095711  0.          ]

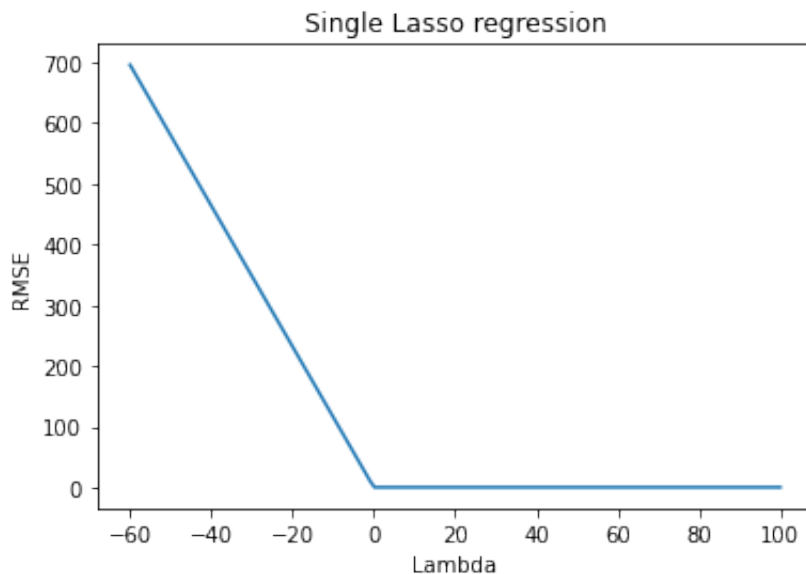
```

▼ Single Lasso Regression

```

1 # Init parameters:
2 xTrain, xTest, yTrain, yTest = train_test_split(np.array(data['yearsofexperie
3 lambdas = np.linspace(-60,100,150)
4
5
6 cont = np.empty([len(lambdas),2])*np.NaN # [lambda error]
7
8 for ii in range(len(lambdas)):
9     lassoModel = Lasso(alpha=lambdas[ii]).fit(xTrain, yTrain)
10    cont[ii,0] = lambdas[ii]
11    error = mean_squared_error(yTest,lassoModel.predict(xTest),squared=False)
12
13    cont[ii,1] = error
14
15 plt.plot(cont[:,0],cont[:,1])
16 plt.xlabel('Lambda')
17 plt.ylabel('RMSE')
18 plt.title('Single Lasso regression')
19 plt.show()
20 optimal_lambda_single = lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))]
21
22 print('Optimal lambda:',lambdas[np.argmax(cont[:,1]==np.min(cont[:,1]))])

```



Optimal lambda: 0.13422818791946156

```

1 #single lasso
2 lassoModel = Lasso(alpha=optimal_lambda_single).fit(xTrain, yTrain)
3 error = mean_squared_error(yTest,lassoModel.predict(xTest),squared=False)
4 r_squared = r2_score(yTest,lassoModel.predict(xTest))
5 print(error,r_squared)
6 print(lassoModel.coef_)

```

```

0.5000023242133684 -0.00010727440064273175
[0.]

```

4. There is controversy as to the existence of a male/female gender pay gap in tech job compensation. Build a logistic

- ▼ regression model (with gender as the outcome variable) to see if there is an appreciable beta associated with total annual compensation with and without controlling for other factors.

▼ Organizing dataset for model and dealing with missing data

Our dataset has both many missing data for the gender column and extremely imbalanced classes. So in this step I am fixing this issue, to input data in the regression that is the least biased as possible.

- I am only using columns that contains an assigned value for gender
- I am randomly sampling all the records with MALE as assigned value for column gender. On this random sample I am using only 6999 records, because it matches the amount available for FEMALE records. This way I am avoiding bias towards a specific class with an imbalanced dataset.

```

1 data.gender = data.gender.str.upper()
2 gender_logit_regression_data = data[(data.gender=='MALE') | (data.gender=='FEM
3

```

```
1 # big difference in the value counts for this category
2 gender_logit_regression_data.gender.value_counts()
```

```
MALE      35698
FEMALE     6999
Name: gender, dtype: int64
```

```
1 # sampling the same amount of records as for female randomly
2 male_df = gender_logit_regression_data[gender_logit_regression_data.gender ==
3
4 female_df = gender_logit_regression_data[gender_logit_regression_data.gender ==
5
6 gender_data = pd.concat([male_df, female_df])
7 gender_data['gender'] = le.fit_transform(gender_data['gender'])
8
9 ## Now Male label stands for 1 and Female label stands for 0
10 gender_data
11
```

	company	title	location	totalyearlycompensation	yearsofexperience
10765	71	11	909	0.036821	0.144928
56603	143	11	58	0.020523	0.173913
56640	640	11	257	0.022133	0.057971
15891	365	8	562	0.048290	0.028986
43268	829	11	824	0.028370	0.072464
...
363	640	11	765	0.033602	0.086957
49636	47	9	625	0.017304	0.086957
27359	47	7	823	0.035211	0.028986
33867	102	1	183	0.018511	0.072464
21922	1000	11	644	0.025352	0.000000

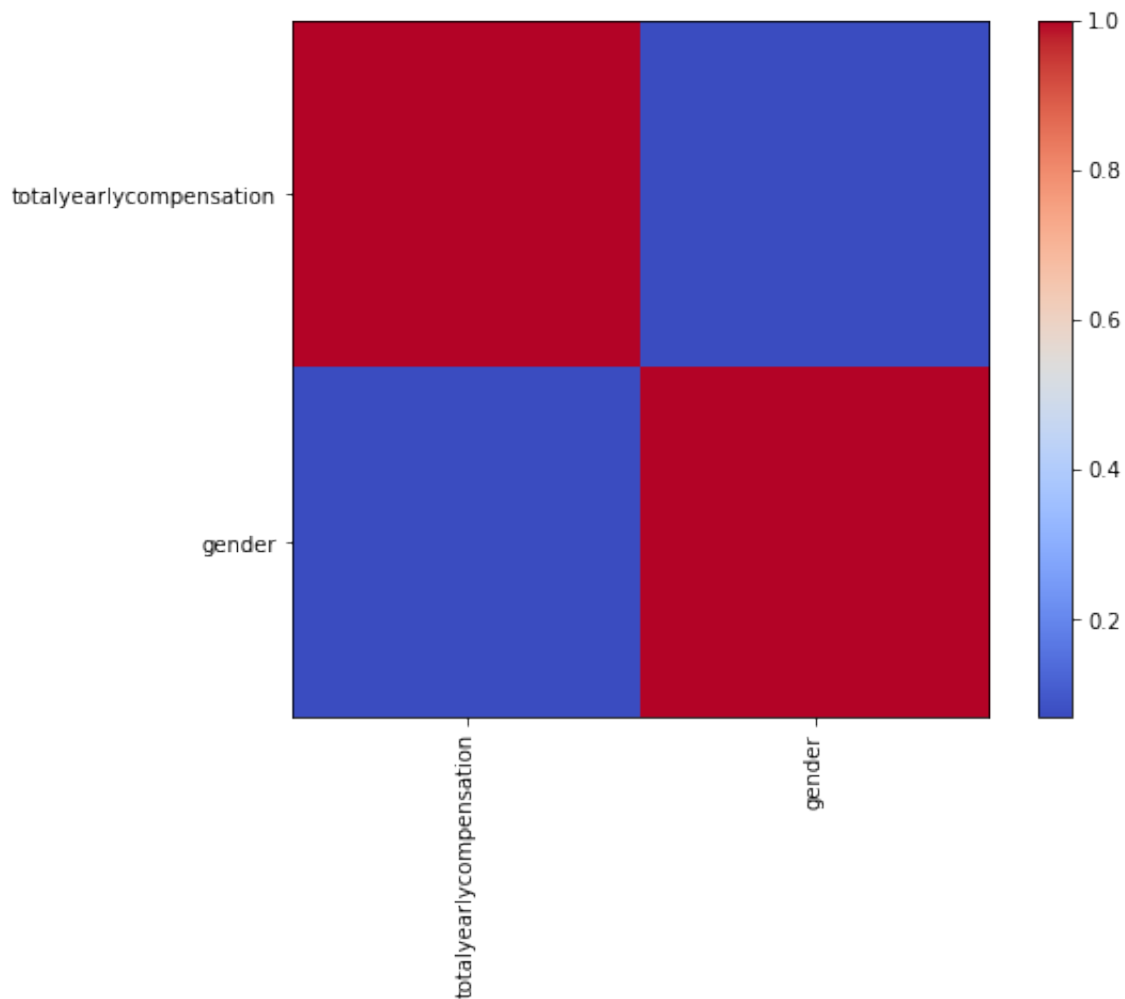
12000 rows x 23 columns

▼ Correlation

```

1 df = pd.DataFrame(gender_data[['totalyearlycompensation','gender']])
2 plt.figure(figsize=(8, 6))# get the correlation matrix
3 corr = df.corr()
4
5 # plot the correlation matrix using matplotlib
6 plt.imshow(corr, cmap='coolwarm', interpolation='nearest')
7 plt.colorbar()
8 plt.xticks(range(len(corr)), corr.columns, rotation=90)
9 plt.yticks(range(len(corr)), corr.columns)
10 plt.show()
11 corr

```



	totalyearlycompensation	gender
totalyearlycompensation	1.000000	0.069288
gender	0.069288	1.000000

▼ The Simple logistical regression

```

1
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import confusion_matrix
5 from sklearn.metrics import matthews_corrcoef
6
7
8 # Split the data into training and testing sets
9 X_train, X_test, y_train, y_test = train_test_split(np.array(gender_data['tot
10
11 # Define your logistic regression model
12 model = LogisticRegression()
13
14 # Train the model on the training data
15 model.fit(X_train, y_train)
16
17 # Test the model on the testing data
18 # predict the dependent variable
19 y_pred = model.predict(X_test)
20
21 # calculate R-squared
22 r_squared = r2_score(y_test, y_pred)
23 rmse = mean_squared_error(y_test, y_pred)
24 acc = model.score(X_test, y_test)
25 cm = confusion_matrix(y_test, y_pred)
26 mcc = matthews_corrcoef(y_test, y_pred)
27
28 print("Accuracy: {:.2f}%".format(acc * 100))
29 print(cm)
30 mcc
31 single_coef = model.coef_[0][0]
32 print(model.coef_)

Accuracy: 52.33%
[[822 341]
 [803 434]]
[[3.60674926]]

```

▼ The multiple logisticat regression

```

1 # Split the data into training and testing sets
2 predictors_columns = ['yearsofexperience', 'title', 'location', 'totalyearlycomp
3 predictors = gender_data.drop(["company", 'title', 'location', 'Race_White', 'Rac
4
5
6 X_train, X_test, y_train, y_test = train_test_split(predictors, np.array(gend
7
8 # Define your logistic regression model
9 model = LogisticRegression()
10
11 # Train the model on the training data
12 model.fit(X_train, y_train)
13 y_pred = model.predict(X_test)
14
15 # Test the model on the testing data
16 r_squared = r2_score(y_test, y_pred)
17 rmse = mean_squared_error(y_test, y_pred)
18 acc = model.score(X_test, y_test)
19 cm = confusion_matrix(y_test, y_pred)
20 mcc = matthews_corrcoef(y_test, y_pred)
21
22 print("Accuracy: {:.2f}%".format(acc * 100))
23 print(cm)
24 mcc
25 print(model.coef_)
26 multiple_coef = model.coef_[0][0]

```

```

Accuracy: 54.00%
[[761 402]
 [702 535]]
[[ 0.0632371  0.32130456  0.07683531 -0.33876056 -0.43016258  0.04844022
    0.09564622  0.05274939  0.00051555 -0.09303457]]

```

```
1 print('change in coefficient', single_coef/multiple_coef)
```

```
change in coefficient 57.035338228713464
```

Interpreting results:

On the first linear regression where we use the `totalyearlycompensation` alone as a predictor for gender. Our model gives us a coefficient of 3.60674926 meaning that for every change on our predictor, the dependable variable will change by a factor of 3.6. This high number must be due to the fact that `totalyearlycompensation` is our only predictor and its doing all of the heavy lifting for the model. We use a correlation, shown on the correlation section of this question, that this variable is very weakly correleated with our target variables.

Subsequently we used a Logistical regression with multiple predictors

[`totalyearlycompensatio`, `yearsofexperience`, `yearsatcompany`, `Masters_Degree`, `Bachelors_Degree`, `Doctorate_Degree`, `Highschool`, `Some_College`, `SAT`, `GPA`].

By comparing the coefficients related to `totalyearlycompensation` we see a decrease by a factor of 57x. Meaning that its 57 times lower than the one on the single logictical regression. While, now our best coeficient is related to the variable `yearsofexperience`.

My interpretation is that `yearsofexperience` might be related to fact that one of the class may have a longer years of experience than the other class. Therefore it works as a better predictor for gender.

Futhermore, we can see that the `totalyearlycompensation` alone cant be used as a single predictor for gender, because once we added extra predictors, our model had a slightly improve.

5. Build a logistic regression model to see if you can predict high

- ▼ and low pay from years of relevant experience, age, height, SAT score and GPA, respectively.

- ▼ Creating cutoff and setting values in the data frame

The cutoff is set by the median value of the `totalyearlycompensation` column. So any value greater than that is beign assigned the class 1 and anythin below is being assigned the class 0.


```

1 pay_Logit_data = data
2
3 def check_cutoff(value):
4     if value > pay_Logit_data.totalyearlycompensation.median():
5         return 1
6     else:
7         return 0
8
9 pay_Logit_data['totalyearlycompensation'] = pay_Logit_data['totalyearlycompensation'].apply(check_cutoff)
10
11 pay_Logit_data

```

	company	title	location	totalyearlycompensation	yearsofexperience
0	722	8	767	0	0.021739
1	331	11	823	0	0.072464
2	47	8	859	1	0.115942
3	71	12	909	1	0.101449
4	640	11	610	0	0.072464
...
62637	437	11	859	1	0.144928
62638	640	11	765	1	0.028986
62639	661	11	859	1	0.202899
62640	832	11	823	1	0.115942
62641	71	11	909	1	0.000000

62637 rows x 23 columns

▼ The logistical regression

By observing the coefficients we can see that as expected from the correlation table in the beginning of the assignment, `yearsofexperience` is one of the most correlated predictors with `totalyearlycompensation`. Therefore it has the highest coefficient. While it is interesting to note that 'Age','Height','GPA' coefficients are negative. Meaning that these predictors have an inverse relation with our target variable.

We can see that with an MCC of 0.38, it is safe to say that our model will be able to predict unseen data, most of the time, and these predictions will be correct.

Accuracy: 69.07%

Coefficients: $\begin{bmatrix} 1.02798357e+01 & -2.08954559e-03 & -2.80691883e-02 & 3.74386274e-03 \\ -1.75117927e-02 \end{bmatrix}$

Confusion Matrix: $\begin{bmatrix} 4616 & 1661 \\ 2214 & 4037 \end{bmatrix}$

MCC: 0.3827647361621342

```
1 predictors_columns = ['yearsofexperience', 'Age', 'Height', 'SAT', 'GPA']
2
3
4 X_train, X_test, y_train, y_test = train_test_split(pay_Logit_data[predictors_
5
6 # Define your logistic regression model
7 model = LogisticRegression()
8
9 # Train the model on the training data
10 model.fit(X_train, y_train)
11 y_pred = model.predict(X_test)
12
13 # Test the model on the testing data
14 acc = model.score(X_test, y_test)
15 cm = confusion_matrix(y_test, y_pred)
16 mcc = matthews_corrcoef(y_test, y_pred)
17 print("Accuracy: {:.2f}%".format(acc * 100))
18 print(model.coef_)
19 print(cm, mcc)
```

```
Accuracy: 69.07%
[[ 1.02798357e+01 -2.08954559e-03 -2.80691883e-02  3.74386274e-03
  -1.75117927e-02]]
[[4616 1661]
 [2214 4037]] 0.3827647361621342
```

Extra credit:

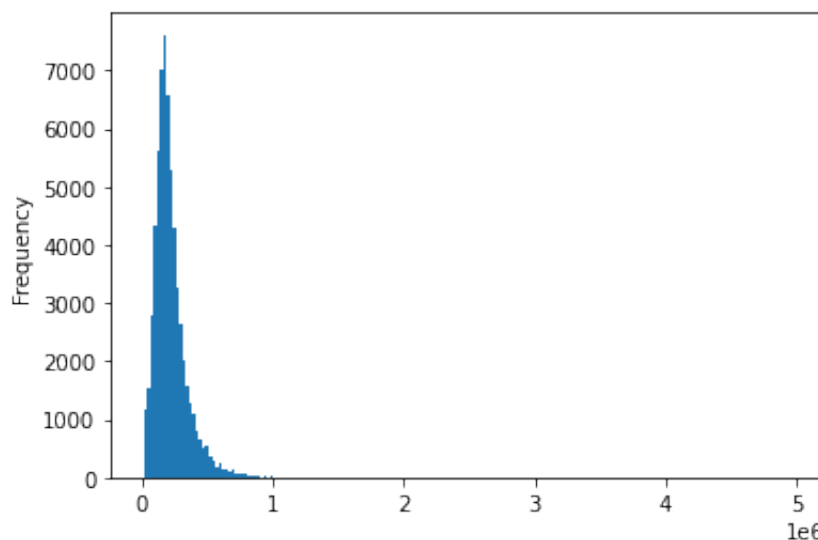
- Is salary, height or age normally distributed? Does this surprise you? Why or why not?
- Tell us something interesting about this dataset that is not already covered by the questions above and that is not obvious.

▼ Extra Credit A

Among all of the 3 variables, height is the only that is normally distributed, as it is a canonical example of the normal distribution given in textbooks. While salary and age are slightly skewed to the left. It is expected that on the case of salary, we will find a large frequency of salaries slightly below the mean, and the higher the salary gets, the lower the frequency. It is expected that age is also skewed in this case, because our data set is about work force. It is expected that the age group of 25 - 40 account for a larger percentage of the worke force. Therefore, there were no surprise given the assumptinos listed above.

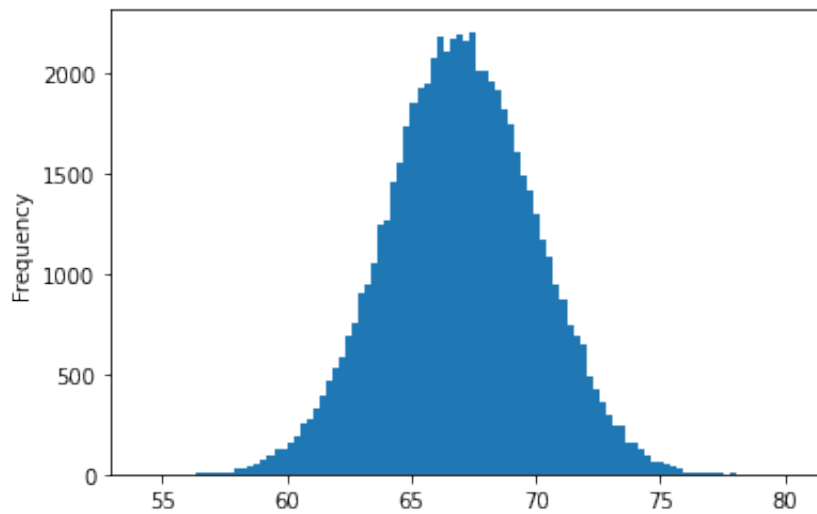
```
1 main_data.totalyearlycompensation.plot.hist(bins=200)
```

<AxesSubplot:ylabel='Frequency'>



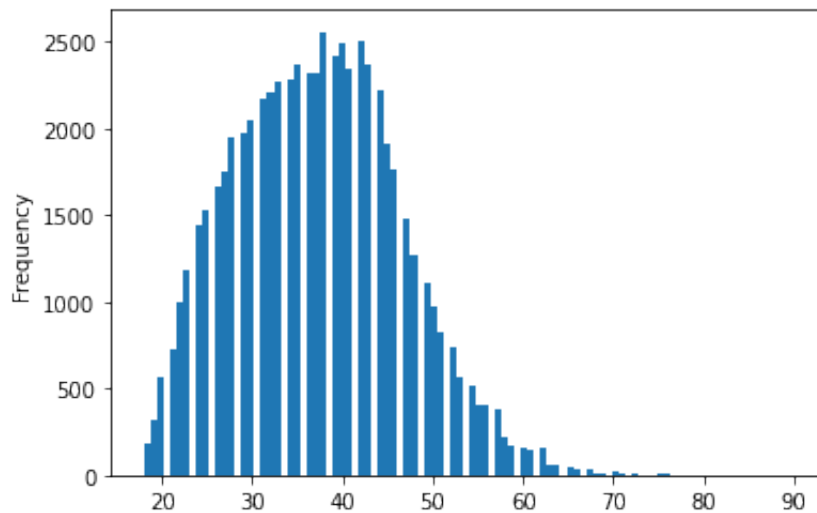
```
1 main_data.Height.plot.hist(bins=100)
```

<AxesSubplot:ylabel='Frequency'>



```
1 main_data.Age.plot.hist(bins=100)
```

<AxesSubplot:ylabel='Frequency'>



▼ Extra Credit B

Something interesting is that if we turn the Title, location and company columns(that are categorical variables) into numeric values, we probably could get a decent model for predicting high and low salaries.

Double-click (or enter) to edit

[Colab paid products](#) - [Cancel contracts here](#)

