

index

January 29, 2022

1 Bias-Variance Tradeoff - Lab

1.1 Introduction

In this lab, you'll practice the concepts you learned in the last lesson, bias-variance tradeoff.

1.2 Objectives

In this lab you will:

- Demonstrate the tradeoff between bias and variance by way of fitting a machine learning model

1.3 Let's get started!

In this lab, you'll try to predict some movie revenues based on certain factors, such as ratings and movie year. Start by running the following cell which imports all the necessary functions and the dataset:

```
[52]: import numpy as np
import pandas as pd

from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import LinearRegression

import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

df = pd.read_excel('movie_data_detailed_with_ols.xlsx')
df.head()
```

```
[52]: Unnamed: 0    budget  domgross    title  Response_Json  Year  \
0          0    13000000  25682380    21 & Over              0   2008
1          1    45658735  13414714    Dredd 3D              0   2012
2          2    20000000  53107035  12 Years a Slave          0   2013
3          3    61000000  75612460    2 Guns              0   2013
```

4	4	40000000	95020213	42	0	2013
---	---	----------	----------	----	---	------

	imdbRating	Metascore	imdbVotes	Model
0	6.8	48	206513	4.912759e+07
1	0.0	0	0	2.267265e+05
2	8.1	96	537525	1.626624e+08
3	6.7	55	173726	7.723381e+07
4	7.5	62	74170	4.151958e+07

Subset the df DataFrame to only keep the 'domgross', 'budget', 'imdbRating', 'Metascore', and 'imdbVotes' columns.

```
[53]: # Subset the DataFrame
to_pick = ['domgross', 'budget', 'imdbRating', 'Metascore', 'imdbVotes']
df = df[to_pick]
df.head()
```

```
[53]:  domgross    budget  imdbRating  Metascore  imdbVotes
0  25682380  13000000         6.8         48      206513
1  13414714   45658735         0.0          0          0
2  53107035  20000000         8.1         96      537525
3  75612460  61000000         6.7         55      173726
4  95020213  40000000         7.5         62       74170
```

1.4 Split the data

- First, assign the predictors to X and the outcome variable, 'domgross' to y
- Split the data into training and test sets. Set the seed to 42 and the test_size to 0.25

```
[54]: # domgross is the outcome variable
X = df.drop('domgross', axis = 1)
y = df['domgross']

X_train , X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size = 0.25,
                                                    random_state=42
                                                    )
```

Use the MinMaxScaler to scale the training set. Remember you can fit and transform in a single method using .fit_transform().

```
[62]: # Transform with MinMaxScaler
scaler = MinMaxScaler()
X_train_scaled = scaler.fit_transform(X_train)
```

Transform the test data (X_test) using the same scaler:

```
[63]: # Scale the test set
X_test_scaled = scaler.fit_transform(X_test)
```

1.5 Fit a regression model to the training data

```
[64]: # Your code
linreg = LinearRegression().fit(X_train_scaled, y_train)
```

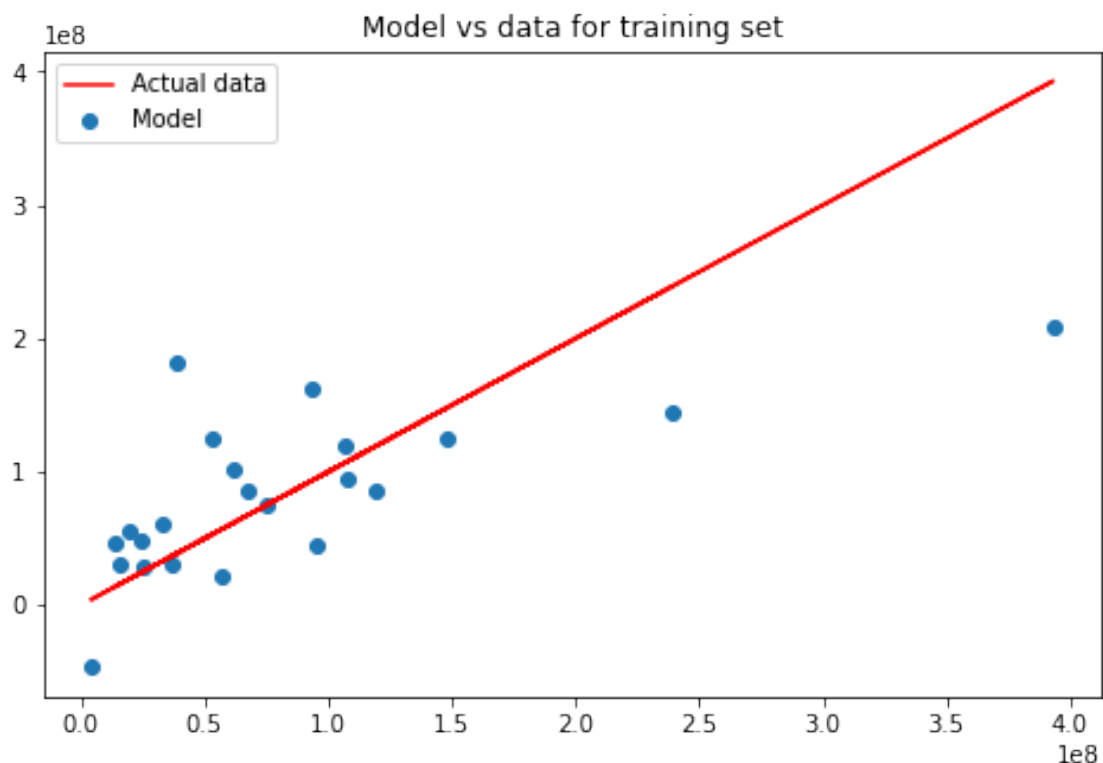
Use the model to make predictions on both the training and test sets:

```
[65]: # Training set predictions
lm_train_predictions = linreg.predict(X_train_scaled)

# Test set predictions
lm_test_predictions = linreg.predict(X_test_scaled)
```

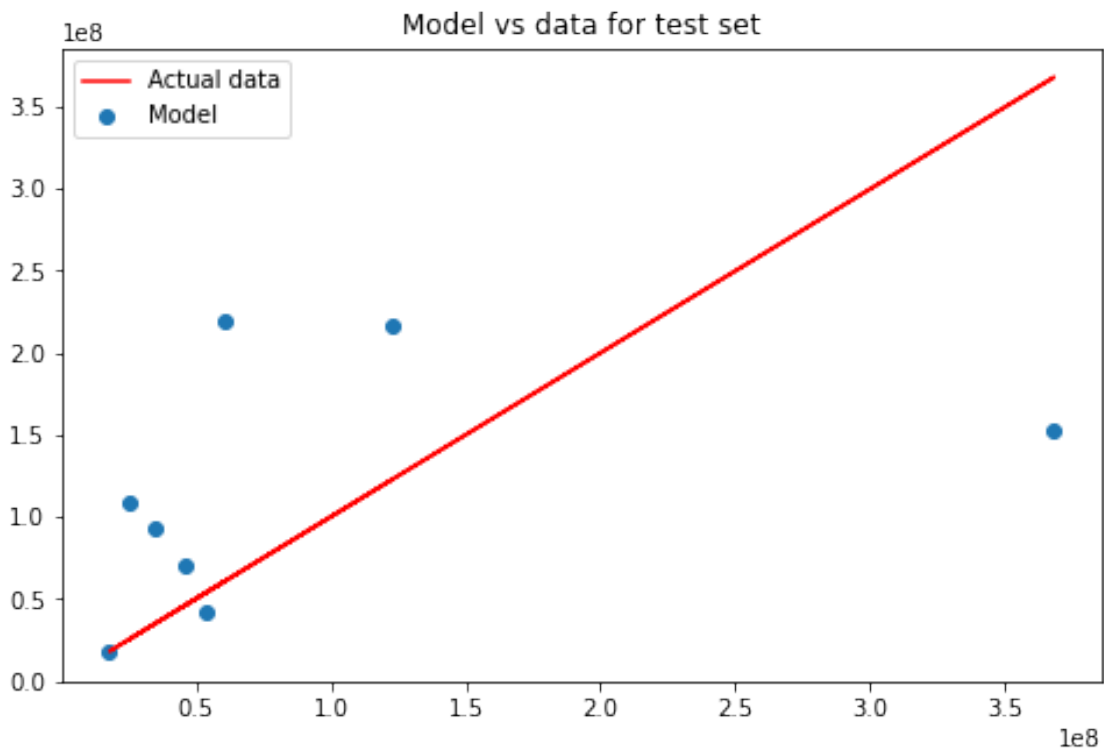
Plot predictions for the training set against the actual data:

```
[66]: # Run this cell - vertical distance between the points and the line denote the ↵
      ↵ errors
plt.figure(figsize=(8, 5))
plt.scatter(y_train, lm_train_predictions, label='Model')
plt.plot(y_train, y_train, label='Actual data', color = "red")
plt.title('Model vs data for training set')
plt.legend();
```



Plot predictions for the test set against the actual data:

```
[67]: # Run this cell - vertical distance between the points and the line denote the errors
plt.figure(figsize=(8, 5))
plt.scatter(y_test, lm_test_predictions, label='Model')
plt.plot(y_test, y_test, label='Actual data', color = "red")
plt.title('Model vs data for test set')
plt.legend();
```



1.6 Bias

Create a function `bias()` to calculate the bias of a model's predictions given the actual data:

$$Bias(\hat{f}(x)) = E[\hat{f}(x) - f(x)] \quad (1)$$

(The expected value can simply be taken as the mean or average value.)

```
[68]: import numpy as np
def bias(y, y_hat):
```

```
return np.mean(y_hat - y)
```

1.7 Variance

Create a function `variance()` to calculate the variance of a model's predictions:

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - (E[\hat{f}(x)])^2 \quad (2)$$

```
[69]: def variance(y_hat):  
       return np.inner(y_hat,y_hat)/len(y_hat) - (np.mean(y_hat))**2
```

1.8 Calculate bias and variance

```
[70]: # Bias and variance for training set  
b = bias(y_train, lm_train_predictions)  
v = variance(lm_train_predictions)  
print('Train bias: {} \nTrain variance: {}'.format(b, v))  
  
# Train bias: -8.127906105735085e-09  
# Train variance: 3406811040986517.0
```

Train bias: -2.7093020352450285e-09

Train variance: 3406811040986519.0

```
[71]: # Bias and variance for test set  
b = bias(y_test, lm_test_predictions)  
v = variance(lm_test_predictions)  
print('Test bias: {} \nTest variance: {}'.format(b, v))  
  
# Test bias: -10982393.918069275  
# Test variance: 1518678846127932.0
```

Test bias: 23766228.24749709

Test variance: 5024352265882432.0

1.9 Overfit a new model

Use `PolynomialFeatures` with degree 3 and transform `X_train_scaled` and `X_test_scaled`.

Important note: By including this, you don't only take polynomials of single variables, but you also combine variables, eg:

\$ Budget * MetaScore ² \$

What you're essentially doing is taking interactions and creating polynomials at the same time! Have a look at how many columns we get using `np.shape()`!

```
[73]: # Your code here
poly = PolynomialFeatures(3)

X_train_poly = poly.fit_transform(X_train)
X_test_poly = poly.fit_transform(X_test)
```

```
[74]: # Check the shape
np.shape(X_train_poly)
```

```
[74]: (22, 35)
```

Fit a regression model to the training data:

```
[75]: # Your code here
polyreg = LinearRegression()
polyreg.fit(X_train_poly, y_train)
```

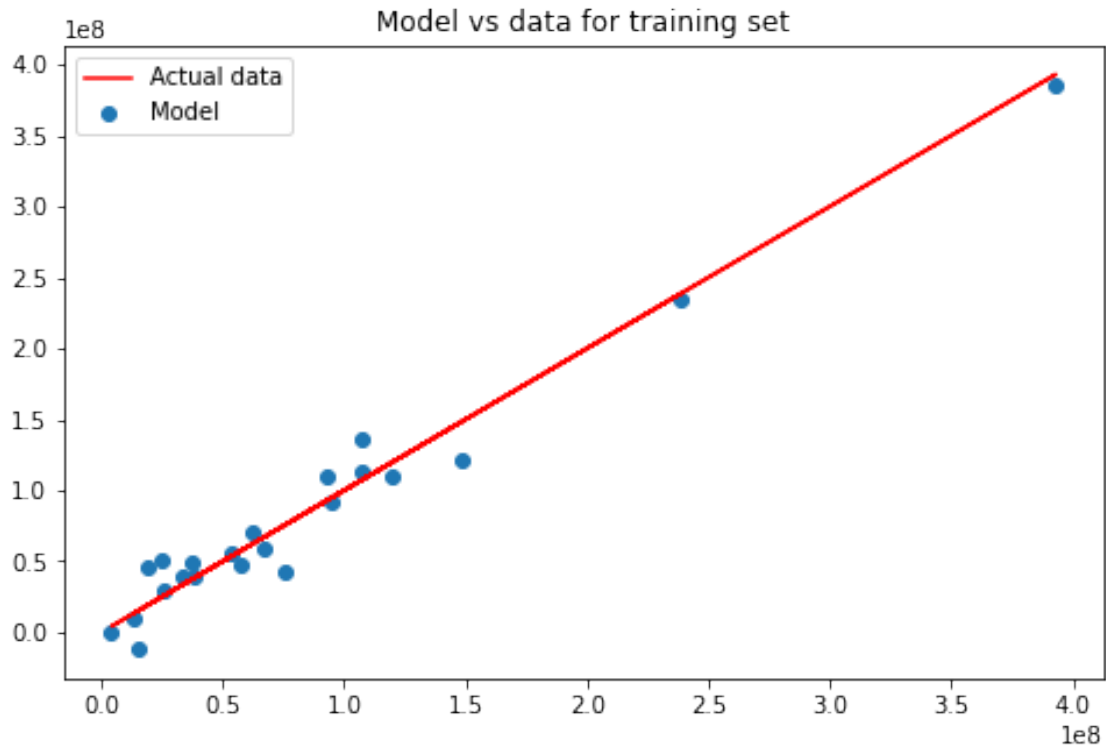
Use the model to make predictions on both the training and test sets:

```
[77]: # Training set predictions
poly_train_predictions = polyreg.predict(X_train_poly)

# Test set predictions
poly_test_predictions = polyreg.predict(X_test_poly)
```

Plot predictions for the training set against the actual data:

```
[79]: # Run this cell - vertical distance between the points and the line denote the
      ↪ errors
plt.figure(figsize=(8, 5))
plt.scatter(y_train, poly_train_predictions, label='Model')
plt.plot(y_train, y_train, label='Actual data', color = "red")
plt.title('Model vs data for training set')
plt.legend();
```



Plot predictions for the test set against the actual data:

```
[80]: # Run this cell - vertical distance between the points and the line denote the
      ↪ errors
      plt.figure(figsize=(8, 5))
      plt.scatter(y_test, poly_test_predictions, label='Model')
      plt.plot(y_test, y_test, label='Actual data', color = "red")
      plt.title('Model vs data for test set')
      plt.legend();
```



Calculate the bias and variance for the training set:

```
[81]: # Bias and variance for training set
b = bias(y_train, poly_train_predictions)
v = variance(poly_train_predictions)
print('Train bias: {} \nTrain variance: {}'.format(b, v))

# Train bias: 3.5898251966996625e-07
# Train variance: 7394168636697528.0
```

Train bias: -1.2598254463889383e-07

Train variance: 7127145906333033.0

Calculate the bias and variance for the test set:

```
[82]: # Bias and variance for test set
b = bias(y_test, poly_test_predictions)
v = variance(poly_test_predictions)
print('Test bias: {} \nTest variance: {}'.format(b, v))

# Test bias: -68166032.47666144
# Test variance: 4.798244829435879e+16
```

Test bias: -66936142.6245157

Test variance: 3.3632039920689276e+16

1.10 Interpret the overfit model

```
[83]: # Your description here

# The training predictions from the second model perfectly match the actual
# data points - which indicates overfitting.

# The bias and variance for the test set both increased drastically for
# this overfit model.
```

1.11 Level Up (Optional)

In this lab we went from 4 predictors to 35 by adding polynomials and interactions, using `PolynomialFeatures`. That being said, where 35 leads to overfitting, there are probably ways to improve by adding just a few polynomials. Feel free to experiment and see how bias and variance improve!

1.12 Summary

This lab gave you insight into how bias and variance change for a training and a test set by using both simple and complex models.