

Polynomial Regression



Animesh Agarwal [Follow](#)

Oct 8, 2018 · 6 min read

This is my third blog in the Machine Learning series. This blog requires prior knowledge of Linear Regression. If you don't know about Linear Regression or need a brush-up, please go through the previous articles in this series.

- [Linear Regression using Python](#)
- [Linear Regression on Boston Housing Dataset](#)

Linear regression requires the relation between the dependent variable and the independent variable to be linear. What if the distribution of the data was more complex as shown in the below figure? Can linear models be used to fit non-linear data? How can we generate a curve that best captures the data as shown below? Well, we will answer these questions in this blog.

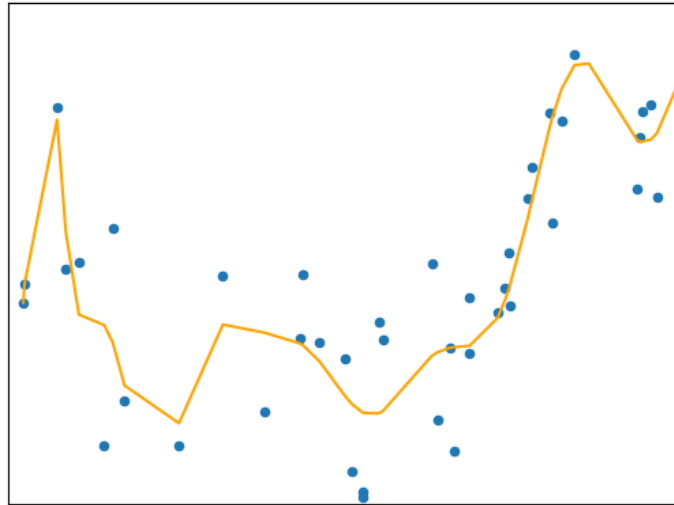


Table of Contents

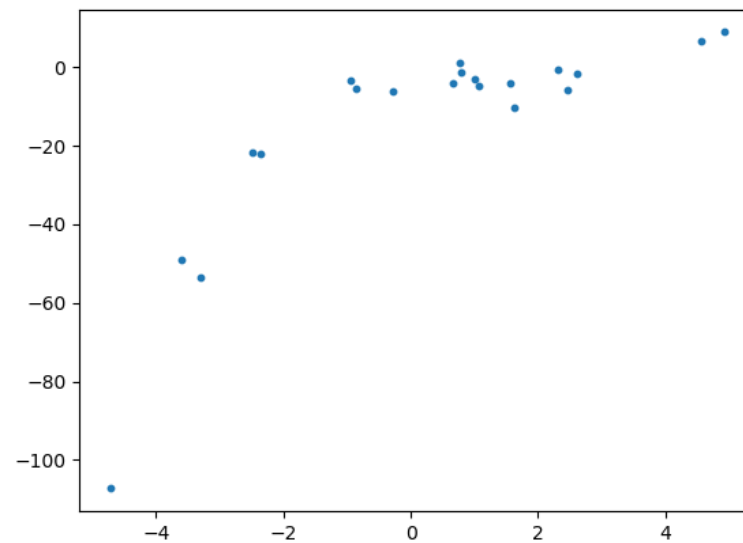
- Why Polynomial Regression
- Over-fitting vs Under-fitting
- Bias vs Variance trade-offs
- Applying polynomial regression to the Boston housing dataset.

Why Polynomial Regression?

To understand the need for polynomial regression, let's generate some random dataset first.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 np.random.seed(0)
5 x = 2 - 3 * np.random.normal(0, 1, 20)
6 y = x - 2 * (x ** 2) + 0.5 * (x ** 3) + np.random.normal(-3,
```

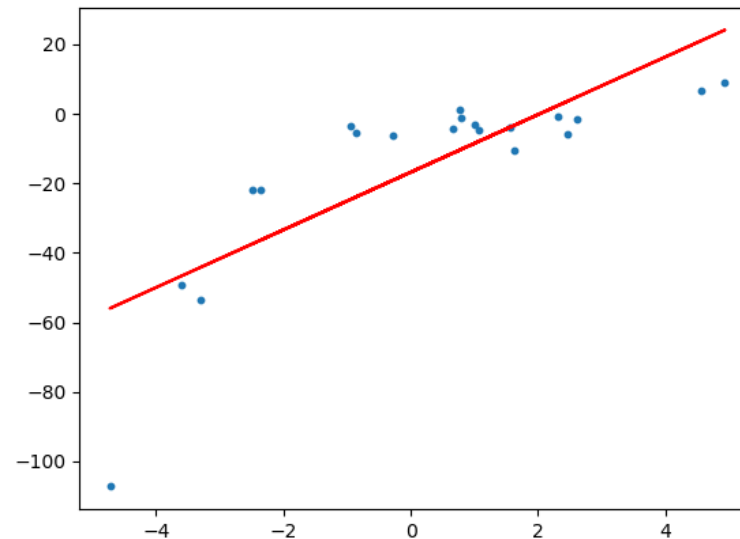
The data generated looks like



Let's apply a linear regression model to this dataset.

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  from sklearn.linear_model import LinearRegression
5
6  np.random.seed(0)
7  x = 2 - 3 * np.random.normal(0, 1, 20)
8  y = x - 2 * (x ** 2) + 0.5 * (x ** 3) + np.random.normal(-3
9
10 # transforming the data to include another axis
11 x = x[:, np.newaxis]
12 y = y[:, np.newaxis]
13
14 model = LinearRegression()
```

The plot of the best fit line is



We can see that the straight line is unable to capture the patterns in the data. This is an example of **under-fitting**. Computing the RMSE and R^2 -score of the linear line gives:

```
RMSE of linear regression is 15.908242501429998.  
R2 score of linear regression is 0.6386750054827146
```

To overcome under-fitting, we need to increase the complexity of the model.

To generate a higher order equation we can add powers of the original features as new features. The linear model,

$$Y = \theta_0 + \theta_1 x$$

can be transformed to

$$Y = \theta_0 + \theta_1 x + \theta_2 x^2$$

*This is still considered to be **linear model** as the coefficients/weights associated with the features are still linear. x^2 is only a feature. However the curve that we are fitting is **quadratic** in nature.*

To convert the original features into their higher order terms we will use the `PolynomialFeatures` class provided by `scikit-learn`. Next, we train the model using Linear Regression.

```
1  import operator
2
3  import numpy as np
4  import matplotlib.pyplot as plt
5
6  from sklearn.linear_model import LinearRegression
7  from sklearn.metrics import mean_squared_error, r2_score
8  from sklearn.preprocessing import PolynomialFeatures
9
10 np.random.seed(0)
11 x = 2 - 3 * np.random.normal(0, 1, 20)
12 y = x - 2 * (x ** 2) + 0.5 * (x ** 3) + np.random.normal(-3
13
14 # transforming the data to include another axis
15 x = x[:, np.newaxis]
16 y = y[:, np.newaxis]
17
18 polynomial_features= PolynomialFeatures(degree=2)
19 x_poly = polynomial_features.fit_transform(x)
20
21 model = LinearRegression()
22 model.fit(x_poly, y)
23 y_poly_pred = model.predict(x_poly)
24
```

To generate polynomial features (here 2nd degree polynomial)

```
-----
polynomial_features = PolynomialFeatures(degree=2)
x_poly = polynomial_features.fit_transform(x)
```

Explanation

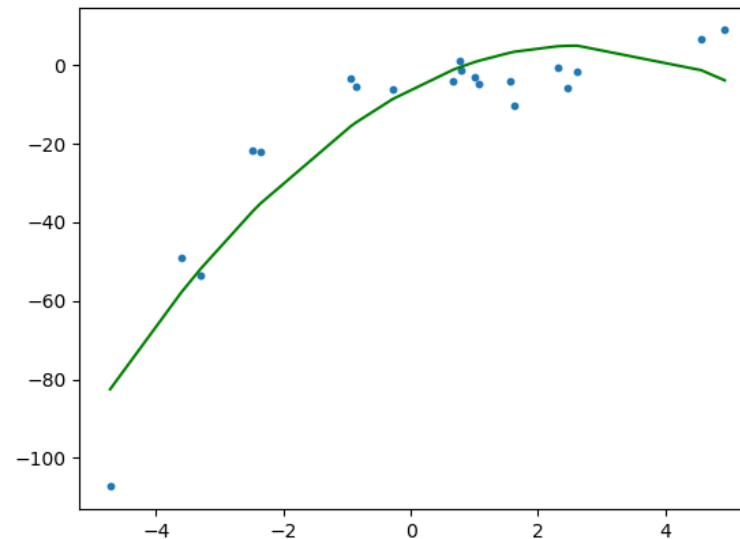
Let's take the first three rows of X:

```
[[-3.29215704]
 [ 0.79952837]
 [-0.93621395]]
```

If we apply polynomial transformation of degree 2, the feature vectors become

```
[[-3.29215704 10.83829796]
 [ 0.79952837  0.63924562]
 [-0.93621395  0.87649656]]
```

Fitting a linear regression model on the transformed features gives the below plot.



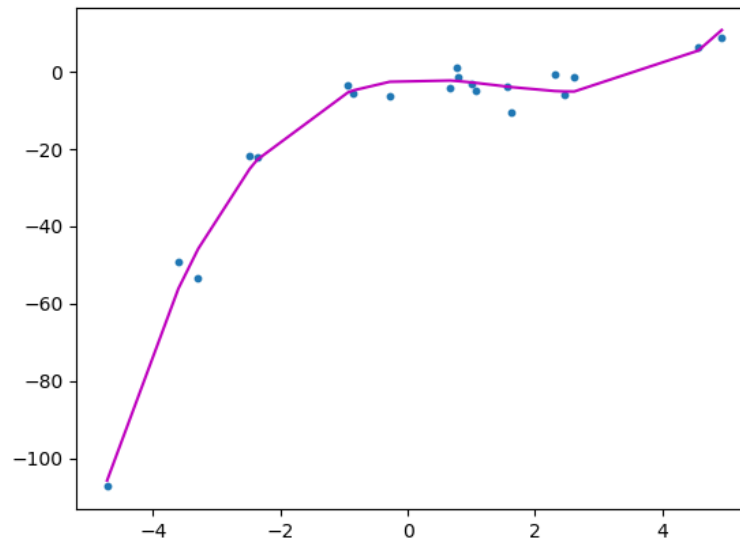
It is quite clear from the plot that the quadratic curve is able to fit the data better than the linear line. Computing the RMSE and R^2 -score of the quadratic plot gives:

```
RMSE of polynomial regression is 10.120437473614711.  
R2 of polynomial regression is 0.8537647164420812.
```

We can see that RMSE has decreased and R^2 -score has increased as compared to the linear line.

If we try to fit a cubic curve (degree=3) to the dataset, we can see that it passes through more data points than the quadratic and the linear

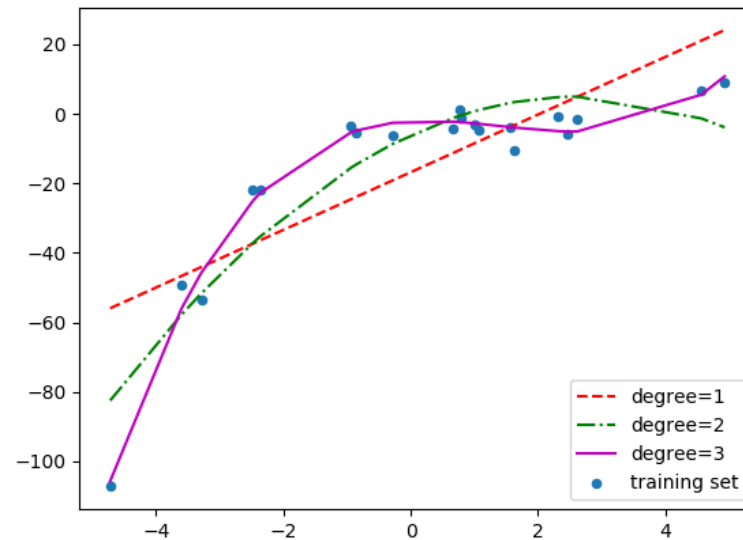
plots.



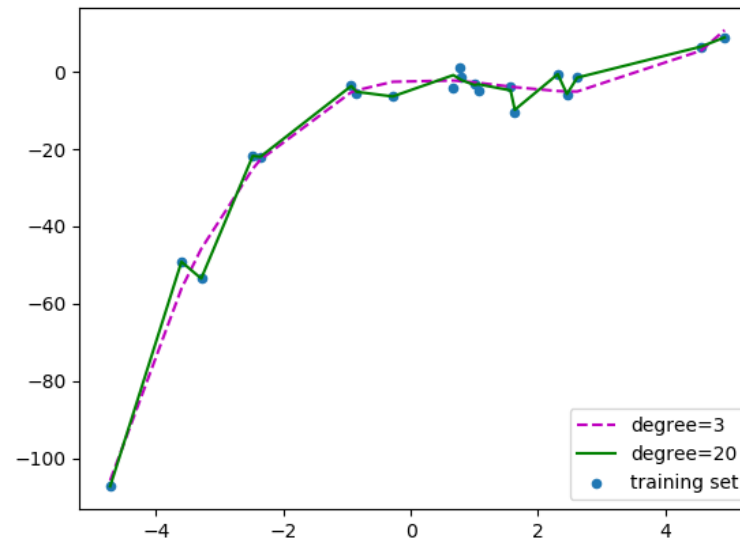
The metrics of the cubic curve is

```
RMSE is 3.449895507408725  
R2 score is 0.9830071790386679
```

Below is a comparison of fitting linear, quadratic and cubic curves on the dataset.



If we further increase the degree to 20, we can see that the curve passes through more data points. Below is a comparison of curves for degree 3 and 20.



For degree=20, the model is also capturing the noise in the data. This is an example of **over-fitting**. Even though this model passes through most of the data, it will fail to generalize on unseen data.

To prevent over-fitting, we can add more training samples so that the algorithm doesn't learn the noise in the system and can become more generalized. (Note: adding more data can be an issue if the data is itself noise).

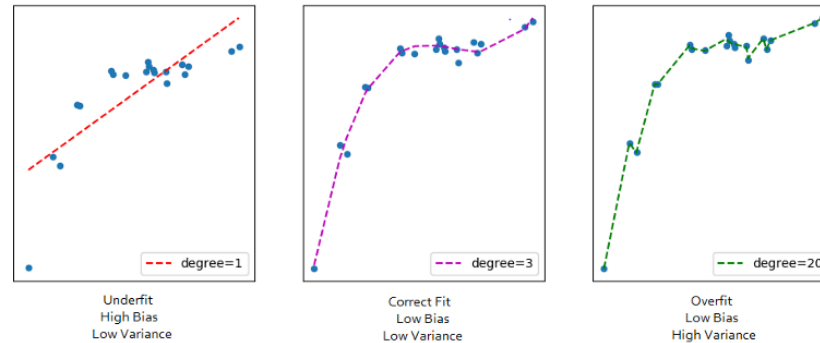
How do we choose an optimal model? To answer this question we need to understand the bias vs variance trade-off.

The Bias vs Variance trade-off

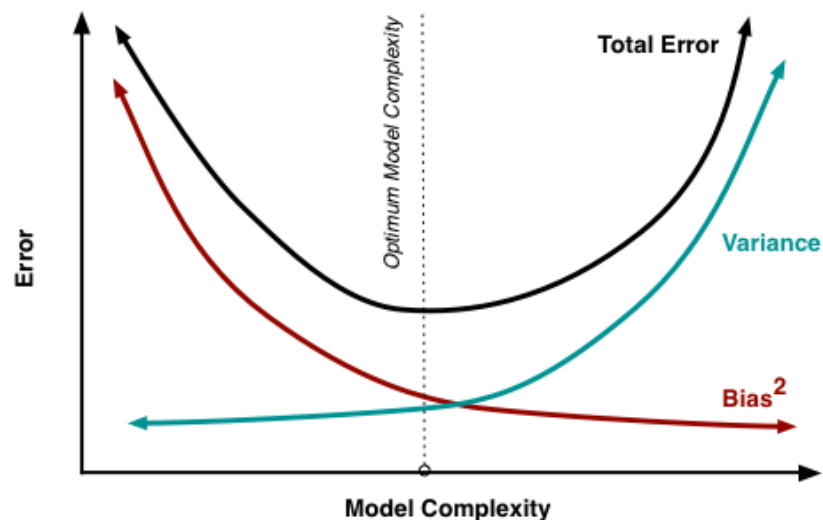
Bias refers to the error due to the model's simplistic assumptions in fitting the data. A high bias means that the model is unable to capture the patterns in the data and this results in **under-fitting**.

Variance refers to the error due to the complex model trying to fit the data. High variance means the model passes through most of the data points and it results in **over-fitting** the data.

The below picture summarizes our learning.



From the below picture we can observe that as the model complexity increases, the bias decreases and the variance increases and vice-versa. Ideally, a machine learning model should have **low variance and low bias**. But practically it's impossible to have both. Therefore to achieve a good model that performs well both on the train and unseen data, a **trade-off** is made.

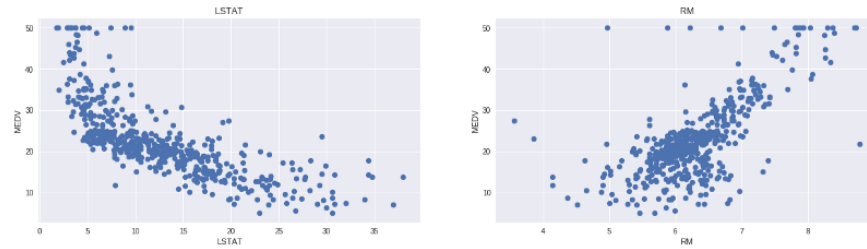


Source: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

Till now, we have covered most of the theory behind Polynomial Regression. Now, let's implement these concepts on the Boston Housing dataset we analyzed in the [previous](#) blog.

Applying Polynomial Regression to the Housing dataset

It can be seen from the below figure that `LSTAT` has a slight non-linear variation with the target variable `MEDV`. We will transform the original features into higher degree polynomials before training the model.



Let's define a function which will transform the original features into polynomial features of a given degree and then apply Linear Regression on it.

```
1  from sklearn.preprocessing import PolynomialFeatures
2
3  def create_polynomial_regression_model(degree):
4      "Creates a polynomial regression model for the given degree"
5
6      poly_features = PolynomialFeatures(degree=degree)
7
8      # transforms the existing features to higher degree features
9      X_train_poly = poly_features.fit_transform(X_train)
10
11     # fit the transformed features to Linear Regression
12     poly_model = LinearRegression()
13     poly_model.fit(X_train_poly, Y_train)
14
15     # predicting on training data-set
16     y_train_predicted = poly_model.predict(X_train_poly)
17
18     # predicting on test data-set
19     y_test_predict = poly_model.predict(poly_features.fit_transform(X_test))
20
21     # evaluating the model on training dataset
22     rmse_train = np.sqrt(mean_squared_error(Y_train, y_train_predicted))
23     r2_train = r2_score(Y_train, y_train_predicted)
24
25     # evaluating the model on test dataset
```

Next, we call the above function with the degree as 2.


```
1 create_polynomial_regression_model(2)
```

The model's performance using Polynomial Regression:

The model performance for the training set

RMSE of training set is 4.703071027847756
R2 score of training set is 0.7425094297364765

The model performance for the test set

RMSE of test set is 3.784819884545044
R2 score of test set is 0.8170372495892174

This is better than what we achieved using Linear Regression in the [previous](#) blog.

That's all for this story. This Github [repo](#) contains all the code for this blog and the complete Jupyter Notebook used for Boston housing dataset can be found [here](#).

Conclusion

In this Machine Learning series, we have covered Linear Regression, Polynomial Regression and implemented both these models on the Boston Housing dataset.

We will cover Logistic Regression in the next blog.

Thanks for Reading !!

