



Bias and Variance in Machine Learning



Renu Khandelwal

Oct 28, 2018 · 6 min read



In this post we will learn how to access a machine learning model's performance.

prerequisites: you need to know basics of machine learning.



First we will understand what defines a model's performance, what is bias and variance, and how bias and variance relate to underfitting and overfitting. Then we will Read more stories this month when you [create a free Medium account](#).

X

How do we know if a model is performing well?

A machine learning model's performance is considered good based on its prediction and how well it generalizes on an independent test dataset. Based on the performance of different models we choose the model which ranks highest in performance.

Let's understand this with an example, let's say we want to predict who will do well in the Midterm election of 2018, will it Republican or Democrats?

We go to a neighborhood and start asking people if they would vote for a Democrat or a Republican. we interview 100 people, 44 say they will vote for Democrats, 40 say they will vote for Republican and 16 are undecided. Based on this data we can make a prediction that chances of Democrats winning is higher than Republicans.

Can we apply this prediction to the entire county, state and then at national level?

No, because the prediction might change if we go to a different neighborhood or a different county or state. We will observe inconsistencies in the prediction. This means our model is not performing well as it cannot be used reliably to make predictions.

One of the reason for our model to perform poorly is due to small sample size and not having enough variation in the data. This introduces error in our prediction. Error is when the predicted value is different from the actual value.

When we have **an input x and we apply a function f on the input x to predict an output y .** Difference between the actual output and predicted output is the error. Our goal with machine learning algorithm is to generate a model which minimizes the error of the test dataset.

Models are assessed based on the prediction error on a new test dataset.

$$L(x, y) = \sum_{i=1}^N (y_i - f(x_i))^2$$

Error = sum of all (Actual output – Predicted Output)

Read more stories this month when you
[create a free Medium account.](#)



Error in our model is summation of reducible and irreducible error.

Error = Reducible Error + Irreducible Error

Reducible Error = Bias² + Variance

Error is ML models known as bias variance decomposition

Irreducible Error

Errors that cannot be reduced no matter what algorithm you apply is called an irreducible error. It is usually caused by unknown variables that may be having an influence on the output variable.

Reducible Error has two components — **bias and variance**.

Presence of bias or variance causes overfitting or underfitting of data.

Bias

Bias is how far are the predicted values from the actual values. If the **average predicted values are far off from the actual values then the bias is high**.

High bias causes algorithm to miss relevant relationship between input and output variable. When a model has a high bias then it implies that the model is too simple and does not capture the complexity of data thus **underfitting the data**.

Variance

Variance occurs when the model performs good on the trained dataset but does not do well on a dataset that it is not trained on, like a test dataset or validation dataset.

Variance tells us how scattered are the predicted value from the actual value.

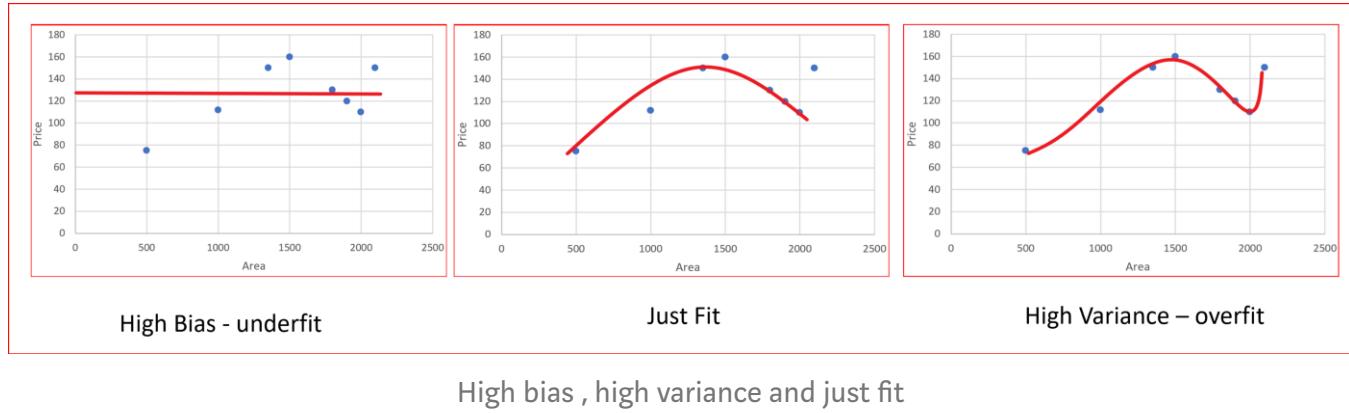
High variance causes overfitting that implies that the algorithm models random noise present in the training data.

when a model has a high variance then the model becomes very flexible and tune itself

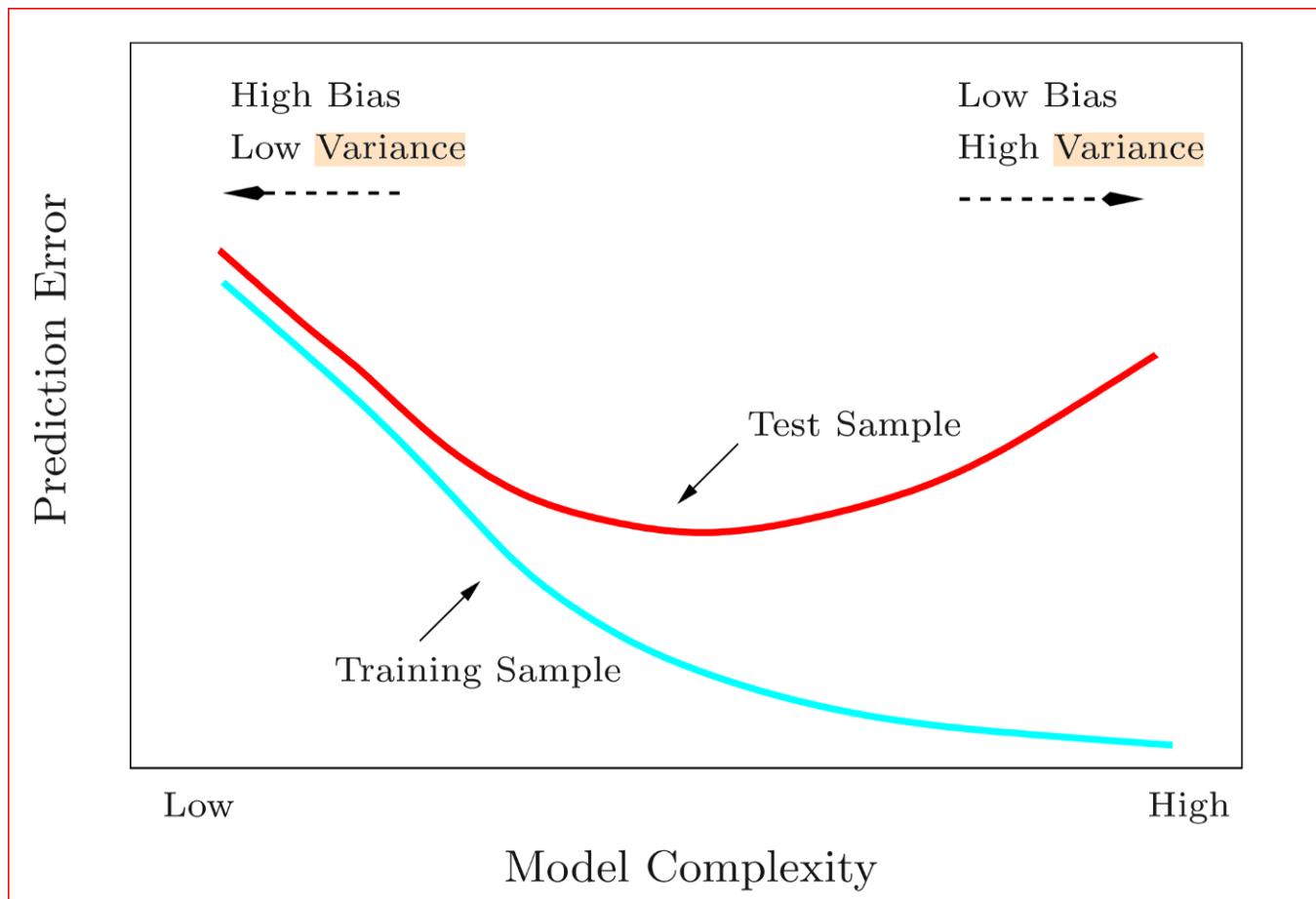
to the data points of the training set when a high variance model overfits to the training set

Read more stories this month when you
[create a free Medium account](#).





If we look at the diagram above, we see that a model with high bias looks very simple. A model with high variance tries to fit most of the data points making the model complex and difficult to model. This can be visible from the plot below between test and training prediction error as a function of model complexity.



Source: Elements of Statistical Learning by Trevor Hastie, Robert Tibshirani and Jerome Friedman

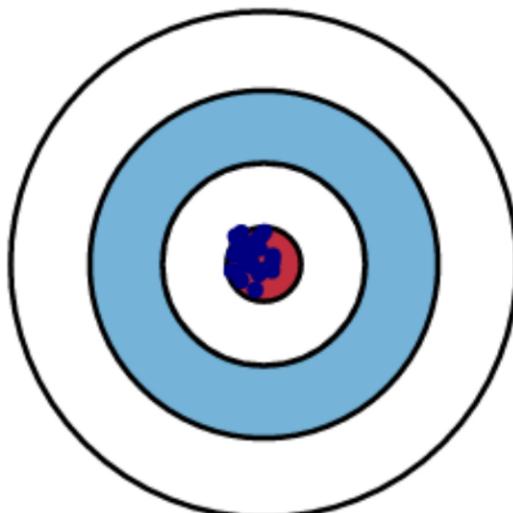
we would like to have a model complexity that trades bias off with variance so that we minimize the test error and would make our model perform better. This is illustrated the the bias variance trade off below.

Read more stories this month when you
[create a free Medium account](#).

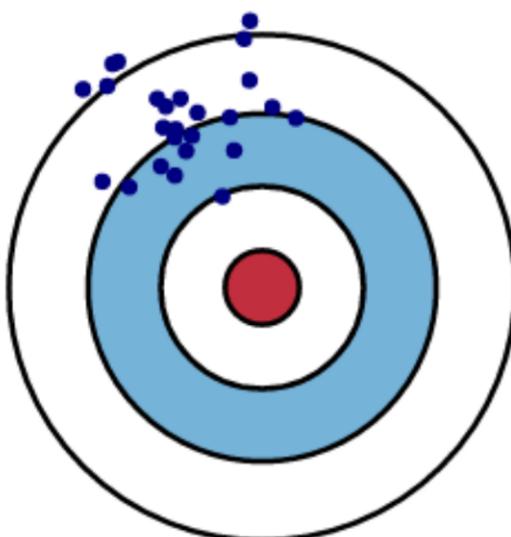
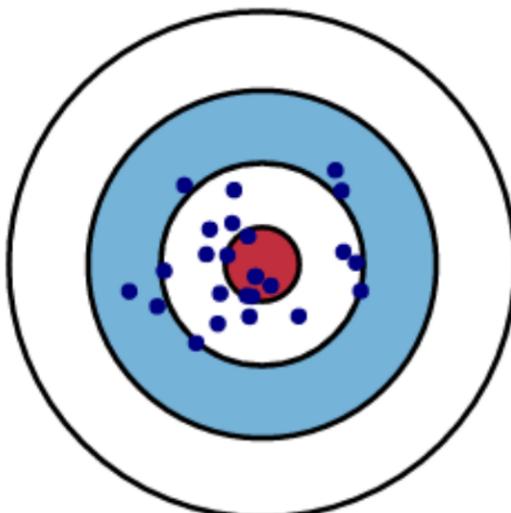
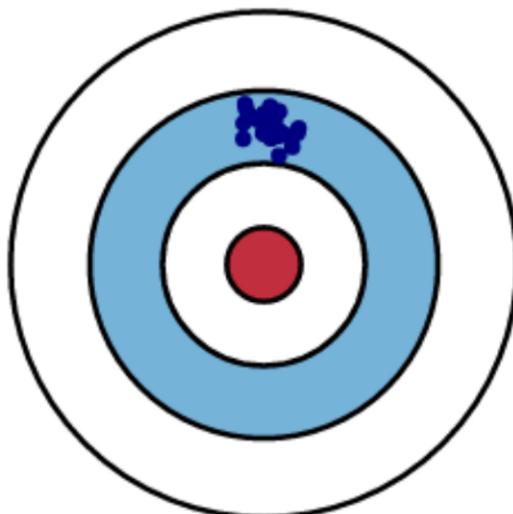


Low Variance High Variance

Low Bias



High Bias



Source: An Introduction to Statistical Learning by Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani

High Bias Low Variance: Models are consistent but inaccurate on average

High Bias High Variance : Models are inaccurate and also inconsistent on average

Low Bias Low Variance: Models are accurate and consistent on averages. We strive for this in our model

Low Bias High variance: Models are somewhat accurate but inconsistent on averages.

Read more stories this month when you
[create a free Medium account.](#)



Is there a way to find when we have a high bias or a high variance?

High Bias can be identified when we have

- High training error
- Validation error or test error is same as training error

High Variance can be identified when

- Low training error
- High validation error or high test error

How do we fix high bias or high variance in the data set?

High bias is due to a simple model and we also see a high training error. To fix that we can do following things

- Add more input features
- Add more complexity by introducing polynomial features
- Decrease Regularization term



Read more stories this month when you
[create a free Medium account.](#)



High bias -Test error is not reduced with more training data

High variance is due to a model that tries to fit most of the training dataset points and hence gets more complex. To resolve high variance issue we need to work on

- Getting more training data
- Reduce input features
- Increase Regularization term



High variance — Test error is reduced with more training data

Before I wind up the topic of bias and variance, a brief on Regularization.

Regularization is a technique where we penalize the loss function for a complex model which is very flexible. This helps with overfitting. It does this by penalizing the different parameters or weights to reduce the noise of the training data and generalize well on the test data

Regularization significantly reduces the variance without substantially increasing bias

Read L1 L2 Regularization here

Read more stories this month when you
[create a free Medium account](#).

×

• • •

Related Posts from DDI:

Deep Learning Explained in 7 Steps - Data Driven Investor

Self-driving cars, Alexa, medical imaging - gadgets are getting super smart around us with the help of deep learning...

www.datadriveninvestor.com

Which is More Promising: Data Science or Software Engineering? - Data Driven Investor

About a month back, while I was sitting at a café and working on developing a website for a client, I found this woman...

www.datadriveninvestor.com



Gain Access to Expert Views

Email

First Name

Give me access!



I agree to leave Medium.com and submit this information, which will be collected and used according to [Upscribe's privacy policy](#).

Read more stories this month when you
[create a free Medium account](#).

×

Machine Learning Bias Variance Bias Variance Trade Off

About Help Legal

Get the Medium app



Read more stories this month when you
[create a free Medium account.](#)

