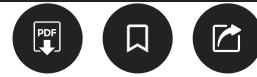


Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning



Deloitte, Flipkart, CRED & other Top Companies are HIRING Data Scientists | 11-13 Feb

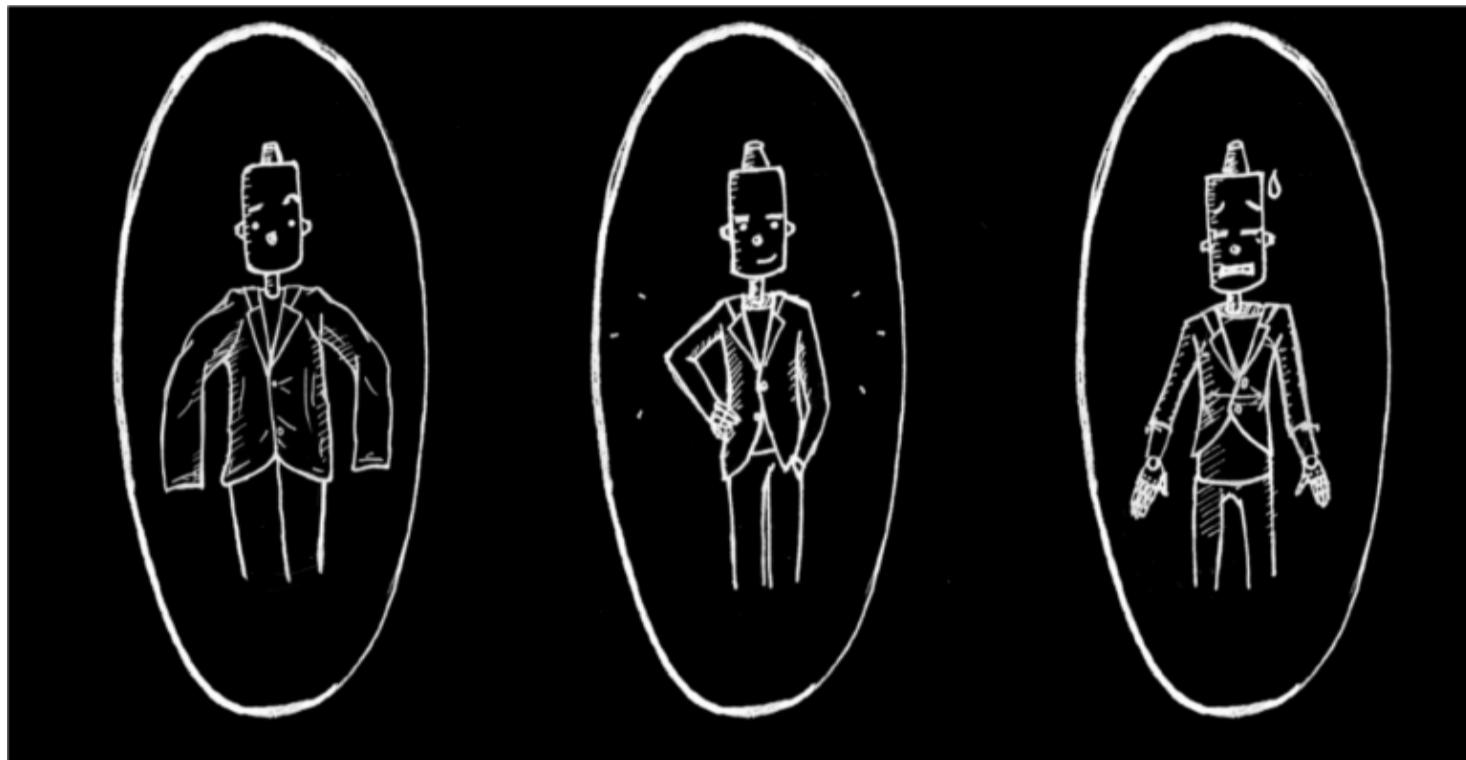
[Register Now](#)



[Home](#)

Sharoon Saxena – February 7, 2020

[Beginner](#) [Machine Learning](#) [Technique](#)



The Challenge of Underfitting and Overfitting in Machine Learning

You'll inevitably face this question in a data scientist interview:

Can you explain what is underfitting and overfitting in the context of machine learning? Describe it in a way even a non-technical person will grasp.

Your ability to explain this in a non-technical and easy-to-understand manner might well decide your fit for the data science role!

Even when we're working on a [machine learning](#) project, we often face situations where we are encountering unexpected performance or error rate differences between the training set and the test set (as shown below). How can a model perform so well over the training set and just as poorly on the test set?

```
1 from sklearn.metrics import classification_report
2 print(classification_report(y_train, predicted_values))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	14234
1	1.00	1.00	1.00	3419
accuracy			1.00	17653
macro avg	1.00	1.00	1.00	17653
weighted avg	1.00	1.00	1.00	17653

```
1 predicted_values = classifier.predict(x_test)
2 print(classification_report(y_test, predicted_values))
```

	precision	recall	f1-score	support
0	0.87	0.86	0.87	3559
1	0.44	0.46	0.45	855
accuracy			0.78	4414
macro avg	0.66	0.66	0.66	4414
weighted avg	0.79	0.78	0.79	4414

This happens very frequently whenever I am working with [tree-based predictive models](#). Because of the way the algorithms work, you can imagine how tricky it is to avoid falling into the overfitting trap!



Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning

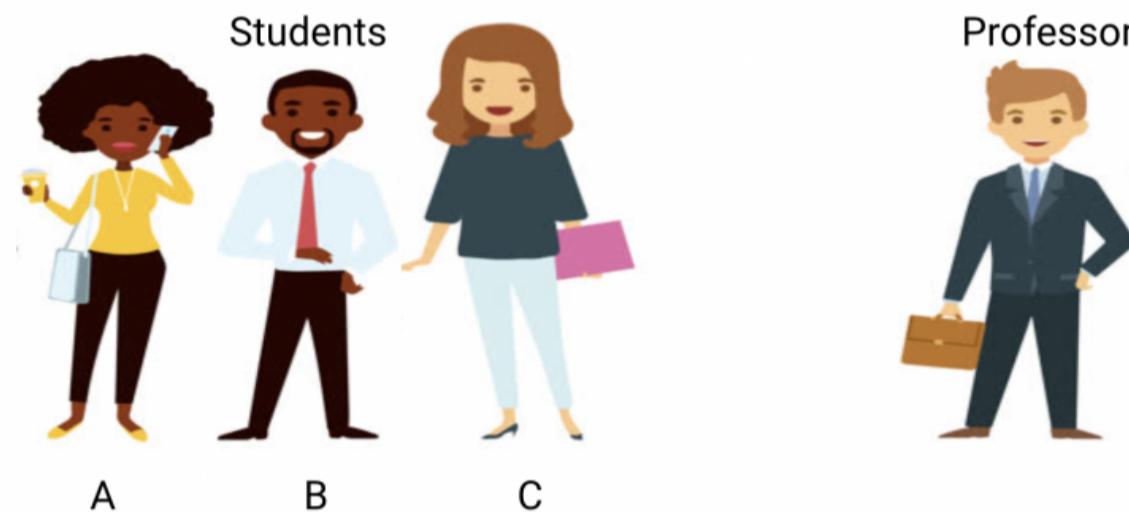
anomalous behavior.

Here's my personal experience – ask any seasoned data scientist about this, they typically start talking about some array of fancy terms like Overfitting, Underfitting, Bias, and Variance. But little does anyone talk about the intuition behind these machine learning concepts. Let's rectify that, shall we?

Let's Take an Example to Understand Underfitting vs. Overfitting

I want to explain these concepts using a real-world example. A lot of folks talk about the theoretical angle but I feel that's not enough – we need to visualize how underfitting and overfitting actually work.

So, let's go back to our college days for this.



Consider a math class consisting of 3 students and a professor.

Now, in any classroom, we can broadly divide the students into 3 categories. We'll talk about them one-by-one.



- Hobby = chating
- Not interested in class
- Doesn't pay much attention to professor

A

Let's say that student A resembles a student who does not like math. She is not interested in what is being taught in the class and therefore does not pay much attention to the professor and the content he is teaching.



- Hobby = to be best in class.
- Mugs up everything professor says.
- Too much attention to the class work.

B

Let's consider student B. He is the most competitive student who focuses on memorizing each and every question being taught in class instead of focusing on the key concepts. Basically, he isn't interested in learning the problem-solving approach.



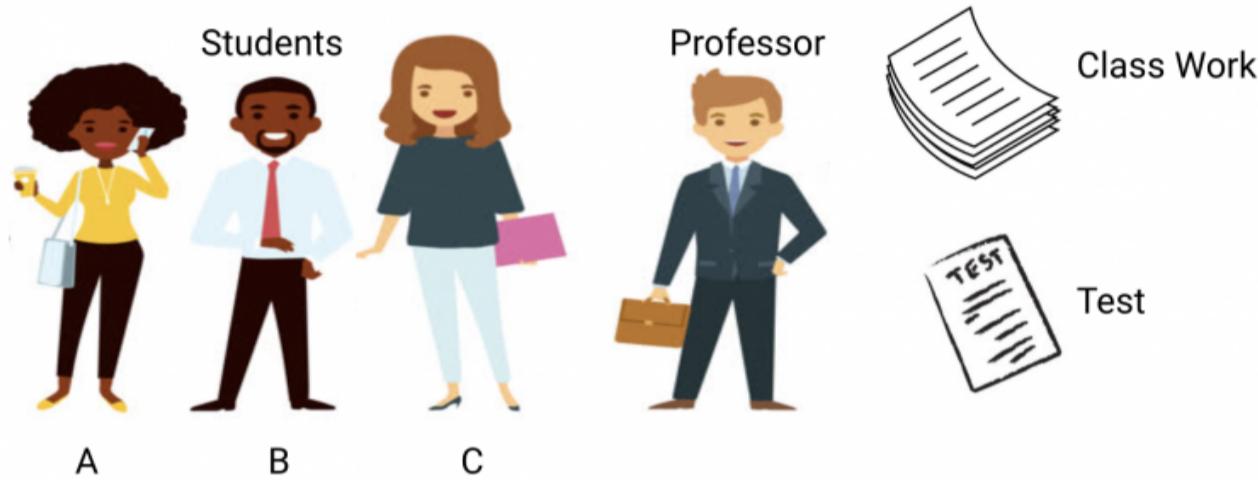
Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning



C

- Eager to learn concepts.
- Pays attention to class and learns the idea behind solving a problem.

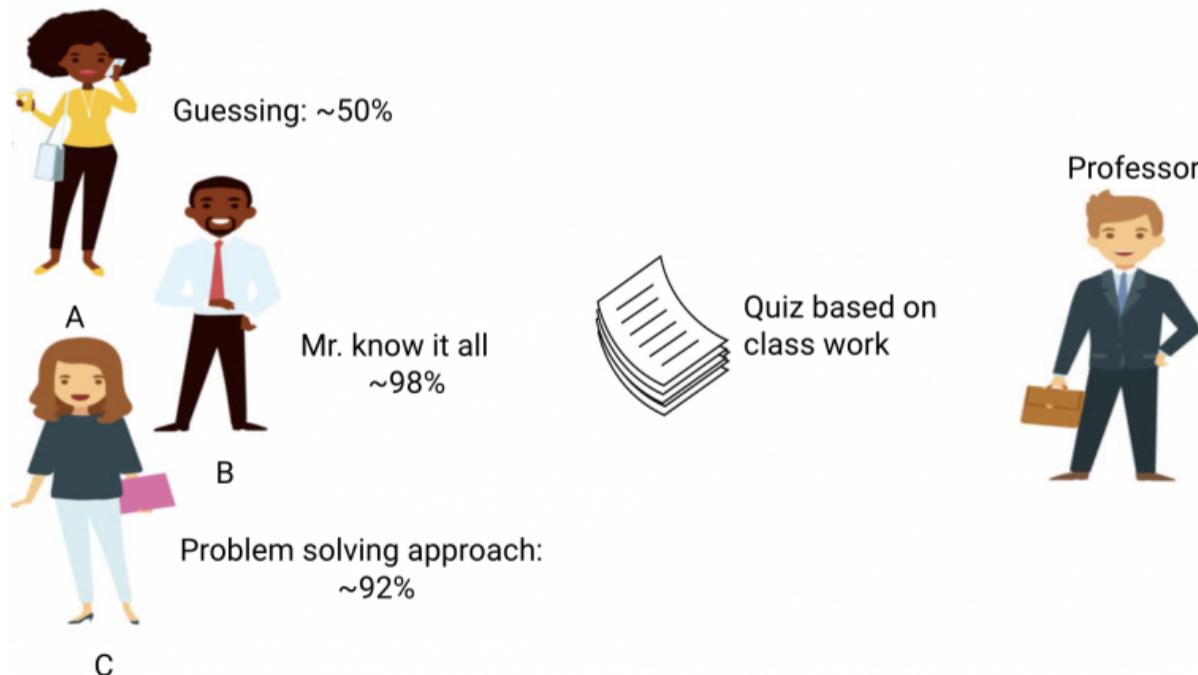
Finally, we have the ideal student C. She is purely interested in learning the key concepts and the problem-solving approach in the math class rather than just memorizing the solutions presented.



We all know from experience what happens in a classroom. The professor first delivers lectures and teaches the students about the problems and how to solve them. At the end of the day, the professor simply takes a quiz based on what he taught in the class.

The obstacle comes in the semester 3 tests that the school lays down. This is where new questions (unseen data) comes up. The students haven't seen these questions before and certainly haven't solved them in the classroom. Sounds familiar?

So, let's discuss what happens when the teacher takes a classroom test at the end of the day:

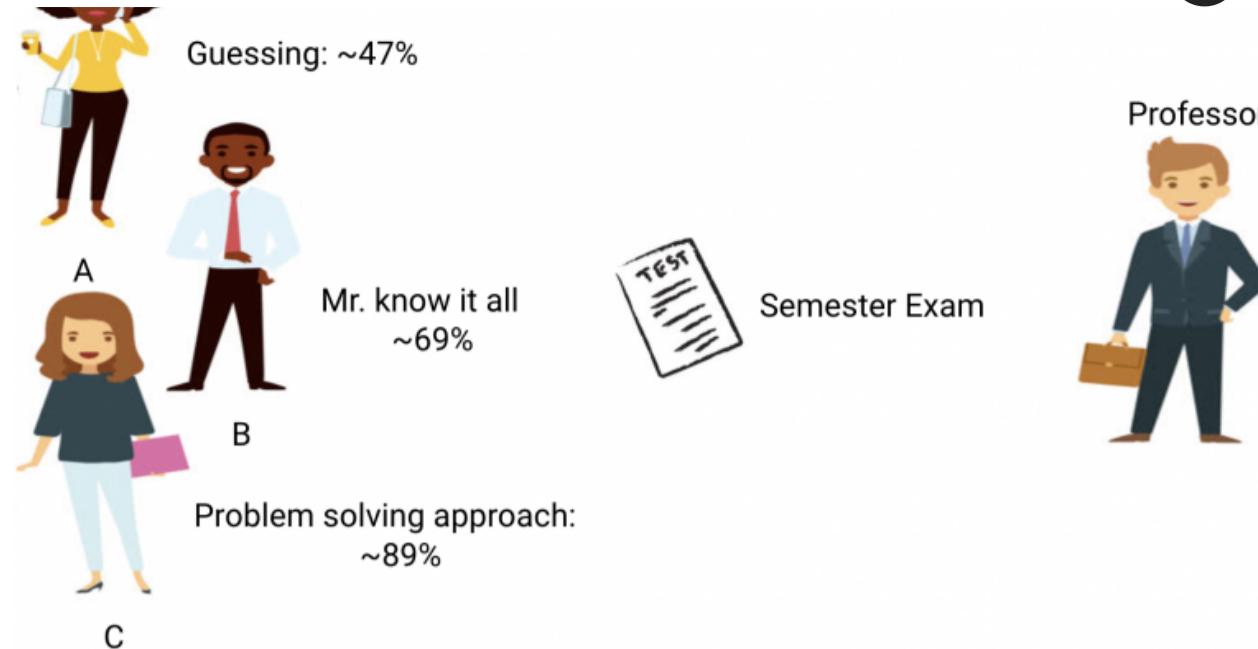
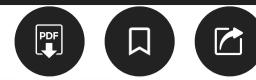


- Student A, who was distracted in his own world, simply guessed the answers and got approximately 50% marks in the test
- On the other hand, the student who memorized each and every question taught in the classroom was able to answer almost every question by memory and therefore obtained 98% marks in the class test
- For student C, she actually solved all the questions using the problem-solving approach she learned in the classroom and scored 92%

We can clearly infer that the student who simply memorizes everything is scoring better without much difficulty.

Now here's the twist. Let's also look at what happens during the monthly test, when students have to face new unknown questions which are not taught in the class by the teacher.

Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning



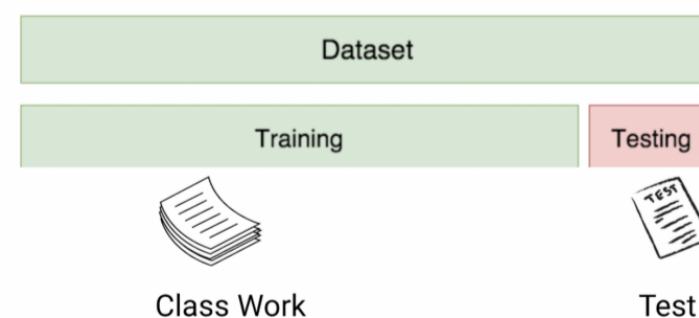
- In the case of student A, things did not change much and he still randomly answers questions correctly ~50% of the time.
- In the case of Student B, his score dropped significantly. Can you guess why? This is because he always memorized the problems that were taught in the class but this monthly test contained questions which he has never seen before. Therefore, his performance went down significantly
- In the case of Student C, the score remained more or less the same. This is because she focused on learning the problem-solving approach and therefore was able to apply the concepts she learned to solve the unknown questions

How Does this Relate to Underfitting and Overfitting in Machine Learning?

You might be wondering how this example relates to the problem which we encountered during the train and test scores of the decision tree classifier? Good question!

	A	B	C
Not interested in learning			
Class test ~50%			
Test ~47%			
Memorizing the lessons			
Class test ~98%			
Test ~69%			
Conceptual Learning			
Class test ~92%			
Test ~89%			

So, let's work on connecting this example with the results of the decision tree classifier that I showed you earlier.



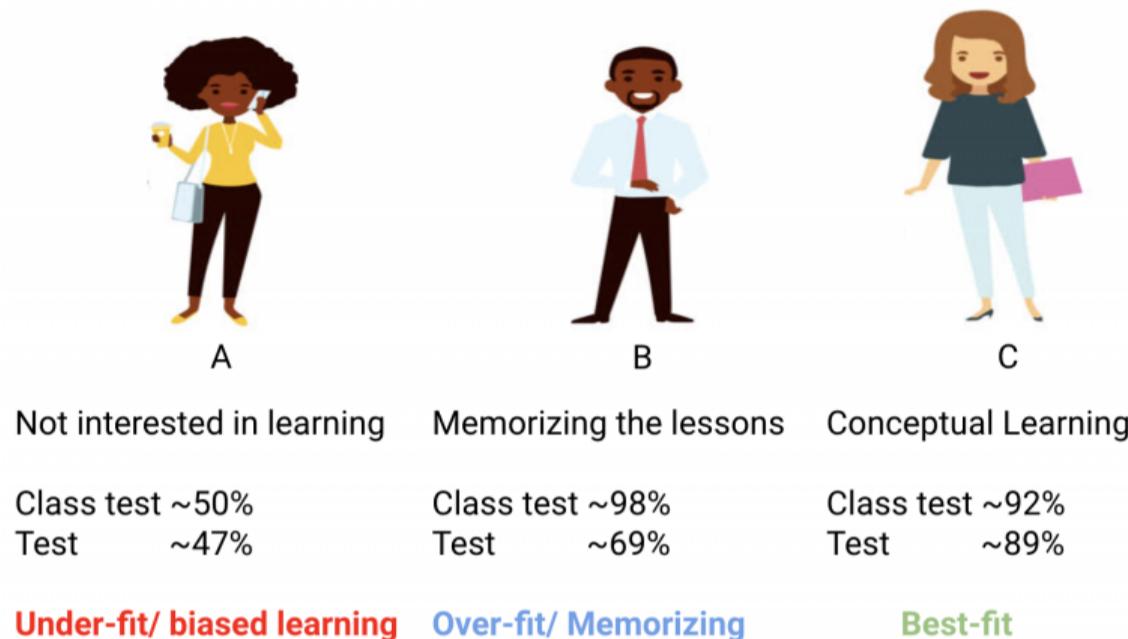
First, the classwork and class test resemble the training data and the prediction over the training data itself respectively. On the other hand, the semester test represents the test set from our data which we keep aside before we train our model (or unseen data in a real-world machine learning project).



Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning

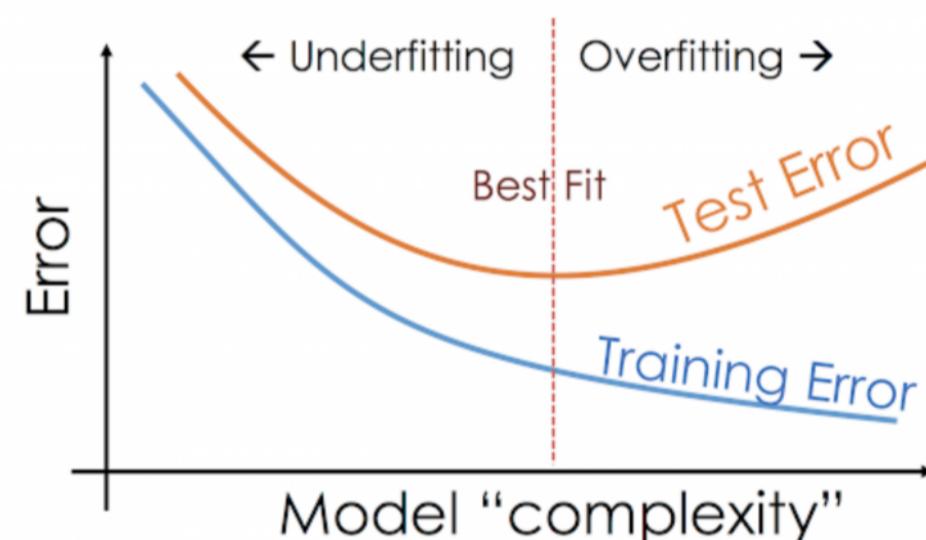
test set. Comparing that to the student examples we just discussed, the classifier establishes an analogy with student B who tried to memorize each and every question in the training set.

Similarly, our decision tree classifier tries to learn each and every point from the training data but suffers radically when it encounters a new data point in the test set. It is not able to generalize it well.



This situation where any given model is performing too well on the training data but the performance drops significantly over the test set is called an overfitting model.

For example, non-parametric models like [decision trees](#), [KNN](#), and [other tree-based algorithms](#) are very prone to overfitting. These models can learn very complex relations which can result in overfitting. The graph below summarises this concept:



On the other hand, if the model is performing poorly over the test and the train set, then we call that an underfitting model. An example of this situation would be building a linear regression model over non-linear data.





Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning

End Notes

I hope this short intuition has cleared up any doubts you might have had with underfitting, overfitting, and best-fitting models and how they work or behave under the hood.

Feel free to shoot me any questions or thoughts below.

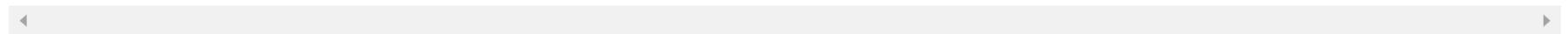
[generalization](#) [overfitting](#) [underfitting](#) [underfitting v overfitting](#)

About the Author



[Sharoon Saxena](#)

Our Top Authors



Download

Analytics Vidhya App for the Latest blog/Article



Previous Post

[Hands-On Tutorial to Analyze Data using Spark SQL](#)

Next Post

[Everything you Need to Know About Scikit-Learn's Latest Update \(with Python Implementation\)](#)

11 thoughts on "Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning"



Thothathiri says:

February 07, 2020 at 11:36 am

I like the way its explained. Well Understood. Kudos to Writer

[Reply](#)



Vered says:

February 07, 2020 at 12:35 pm

Great explanation, thanks! I believe u have a minor mistake in the third quote - it should be "... if the model is performing poorly..."

[Reply](#)



Shiva Golla says:

February 07, 2020 at 8:14 pm



Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning



Manish Thaker says:

February 08, 2020 at 9:59 am

Simply Superb Analogy Sharoon. Expecting more such articles from you. Keep it up.

[Reply](#)



Rashi says:

February 08, 2020 at 10:03 am

Kudos to the explanation

[Reply](#)



Sharoon Saxena says:

February 08, 2020 at 1:06 pm

As I covered in the article, the underfitting and overfitting can be identified using a test set or a validation set from the data. We first train the model on training set and record the performance. Next we also generate predictions over the test set and look over the performance. If Test <<< Train ~ Overfitting If Test~Train and both scores are significantly high in magnitude ~ Good Fit if Test~Train and both scores are below par - Underfit

[Reply](#)



Sharoon Saxena says:

February 08, 2020 at 1:07 pm

thank you for pointing that out, corrected!

[Reply](#)



Ryan Tabeshi says:

February 13, 2020 at 11:57 am

Very well written

[Reply](#)



Yaser Sakkaf says:

February 19, 2020 at 11:05 am

Good one. Well Explained.

[Reply](#)



Rajendra Singh Nayal says:

June 24, 2020 at 12:08 pm

Good & intuitive way to explain over-fitting and under-fitting. I however didn't like the use of He/ his for student 'A' who has been shown as a girl. I know this doesn't matter for the purpose of the article but still it will be nice if this issue can be fixed.

[Reply](#)



Ajay Kopperla says:

August 30, 2020 at 7:32 pm

Very Nice Explanation. It explained the concept in very simple terms

[Reply](#)

Leave a Reply

Your email address will not be published. Required fields are marked *

Comment

Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning

[Website](#)

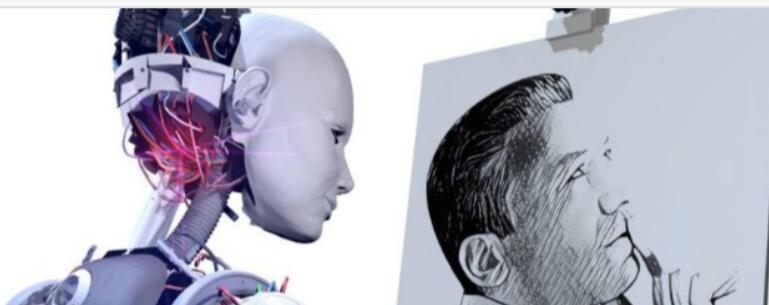


Notify me of follow-up comments by email.

Notify me of new posts by email.

[Submit](#)

Top Resources



[An Introduction to Synthetic Image Generation from Text Data](#)

Suvojit Hore - JAN 28, 2022



[Python Tutorial: Working with CSV file for Data Science](#)

Harika Bonthu - AUG 21, 2021



[3 Interesting Python Projects With Code for Beginners!](#)

Gaurav Sharma - JUL 18, 2021



[Commonly used Machine Learning Algorithms \(with Python and R Codes\).](#)

Sunil Ray - SEP 09, 2017



Download App



Analytics Vidhya

[About Us](#)

[Our Team](#)

[Careers](#)

[Contact us](#)

Companies

[Post Jobs](#)

[Trainings](#)

[Hiring Hackathons](#)

[Advertising](#)

Data Scientists

[Blog](#)

[Hackathon](#)

[Discussions](#)

[Apply Jobs](#)

[Visit us](#)



Underfitting vs. Overfitting (vs. Best Fitting) in Machine Learning

