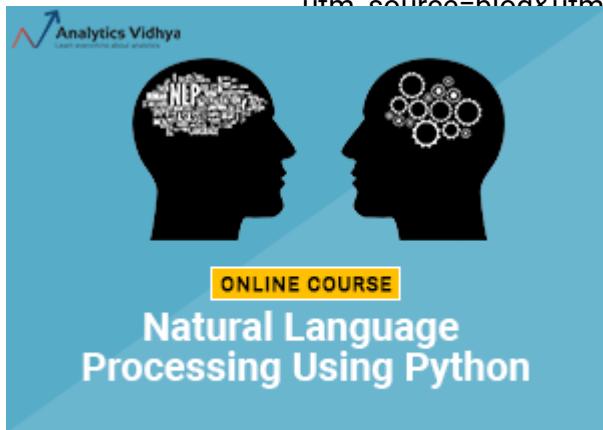


[FLAT 25% OFF on ALL Courses & Programs \(Use Coupon: LETSLEARN\)](#) | [Enroll Today](#)  
([https://courses.analyticsvidhya.com/collections/?utm\\_source=blog&utm\\_medium=flash\\_strip&utm\\_campaign=25\\_off](https://courses.analyticsvidhya.com/collections/?utm_source=blog&utm_medium=flash_strip&utm_campaign=25_off))



 LOGIN / REGISTER ([HTTPS://ID.ANALYTICSVIDHYA.COM/AUTH/LOGIN/?](https://id.analyticsvidhya.com/auth/login/?)

G&UTM MEDIUM=6-PYTHON-LIBRARIES-INTERPRET-MACHINE-LEARNING-MODELS

Learn to Solve  
Text Classification Problems Using **NLP**

<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>



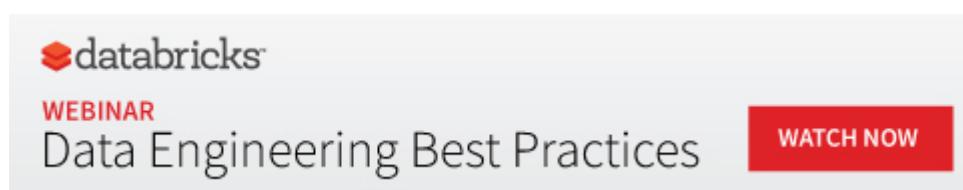
isplay&utm\_campaign=NLPcourse)

(<https://www.analyticsvidhya.com/blog/>)

CATEGORY/INTERMEDIATE/)  
BLOG/CATEGORY/MACHINE-  
GORY/PYTHON-2/)  
TEGORY/TECHNIQUE/)

# Decoding the Black Box: An Important Introduction to Interpretable Machine Learning Models in Python

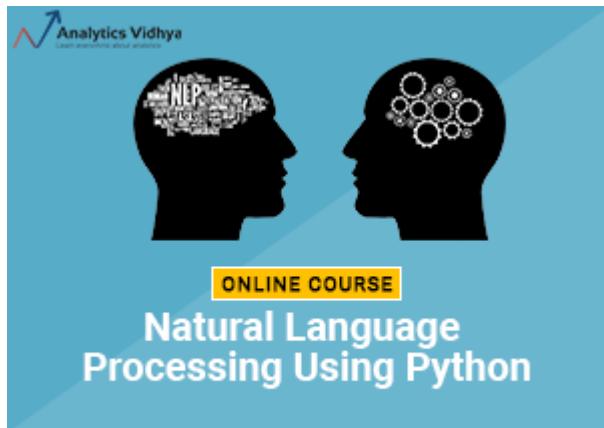
[utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=bain\\_hack](#))  
ANKIT CHOUDHARY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/AUTHOR/ANKIT2106/](https://www.analyticsvidhya.com/blog/author/ankit2106/)), AUGUST 26, 2019 [LOGIN TO BOOKMARK THIS ...](#)



## Overview

- Interpretable machine learning is a critical concept every data scientist should be aware of

- How can you build interpretable machine learning models? This article will provide a framework
- We will also code these interpretable machine learning models in Python



**Learn to Solve  
Text Classification Problems Using NLP**  
[\(https://holerainbowanalytica.com/starter-guide-to-natural-language-processing-for-stakeholders/\)](https://holerainbowanalytica.com/starter-guide-to-natural-language-processing-for-stakeholders/)



Can you imagine building a facial recognition software that misclassifies a person? Or a credit card fraud detection model that raises an alarm for a perfectly legal transaction? And then not being able to explain why that's happening – not ideal.

So the question is – how do we build interpretable machine learning models? That's what we will talk about in this article. We'll first understand what interpretable machine learning is and why it's important. Then we will understand a simple framework for interpretable ML and use that to build machine learning models.



[https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm\\_source=blockutm\\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python](https://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?utm_source=blockutm_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python). The course provides you all the tools and techniques you need to solve business problems using NLP for beginners as well as intermediate-level professionals!



[https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=bain\\_hack](https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)

- Model Agnostic Techniques for Interpretable Machine Learning
- LIME (Local Interpretable Model Agnostic Explanations)
- Python Implementation of Interpretable Machine Learning Techniques

## What is Interpretable Machine Learning?

How do we build trust in machine learning models? That's essentially what this boils down to.

Machine learning-powered applications have become an ever-increasing part of our lives, from image and facial recognition systems to conversational applications, autonomous machines, and personalized systems.

**Analytics Vidhya**  
Learn everything about analytics

**ONLINE COURSE**

## Natural Language Processing Using Python

Learn to Solve Text Classification Problems Using **NLP**

<https://www.analyticsvidhya.com/courses/natural-language-processing-in-python>

Instantly Convert EDW Big Data/Cloud  
Automatic Conversion vs. Traditional Approach

**CONTACT US**

**IMPETUS**

**Women-in-the-loop**  
A Data Science Hackathon by Bain & Company

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

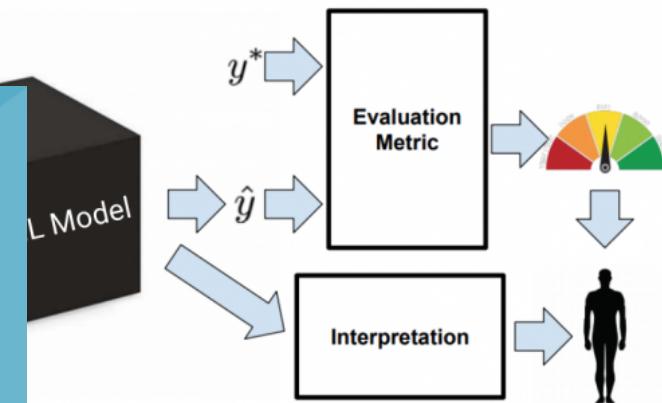
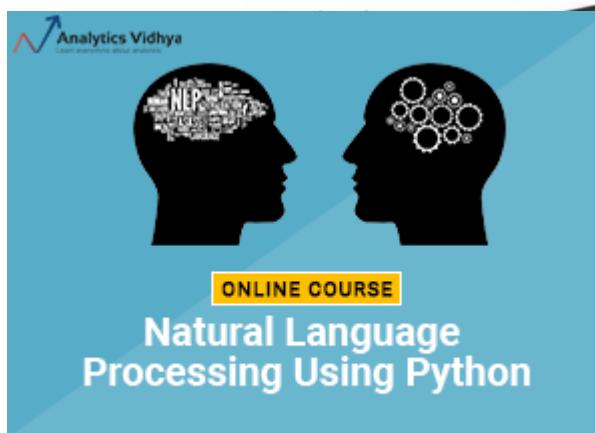
Items based on structured datasets. However, a key issue that arises especially as model complexity increases is interpretability. Let's start with the formal definition:

in-the-loop-a-data-science-hackathon-by-

bain/?

utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

*Interpretation of a machine learning model is the process wherein we try to understand the predictions of a machine learning model.*



[/uploads/2019/08/Screenshot-2019-08-16-at-3.07.26-PM.png\).](#)

Learn to Solve  
Text Classification Problems Using **NLP**  
(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

Modeling lifecycle is at two important stages:

One – where we monitor the evaluation metric and try different ideas of feature engineering, feature selection and build more robust models  
The other – where we interpret the models using the predictions and classifier chose a particular class for example



## What is Machine Learning?

Machine learning is a process for interpreting machine learning models  
[https://www.analyticsvidhya.com/blog/2019/08/applied-machine-learning-beginner-to-professional/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=bain\\_hack](https://www.analyticsvidhya.com/blog/2019/08/applied-machine-learning-beginner-to-professional/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)  
[https://www.analyticsvidhya.com/blog/2019/08/black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm\\_source=...](https://www.analyticsvidhya.com/blog/2019/08/black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm_source=...)

Let's look at why this is important  
<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>

**Relive your favorite memories.**

Nest Hub



## Fairness

Let's take a simple example to understand this. Suppose we are trying to predict employees' performance in a big company to expedite the appraisal process and identify the best employees.

We have data from the last 10 years about the performance reviews of employees. But what if that company tends to promote more men than women?

[Learn to Solve Text Classification Problems Using NLP](http://courses.analyticsvidhya.com/courses/natural-language-processing/text-classification-problems-using-nlp) (<http://courses.analyticsvidhya.com/courses/natural-language-processing/text-classification-problems-using-nlp/uploads/2019/08/Screenshot-2019-08-16-at-3.11.25-PM.png>).

language-processing-nlp/?utm\_source=medium&utm\_medium=referral&utm\_campaign=NLPcourse). Now, if there is no way to interpret our model at this point, we can still gain some insights at the cost of compromising on fairness.

<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by/>

bain/? Let's consider another example. Consider that we are building a model for classifying wolves vs dogs. The data (utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack) that is available is simply labeled images of dogs and wolves.



Dog

(<https://coders.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

entirely possible as wolves are while dogs are in completely different backgrounds (households)

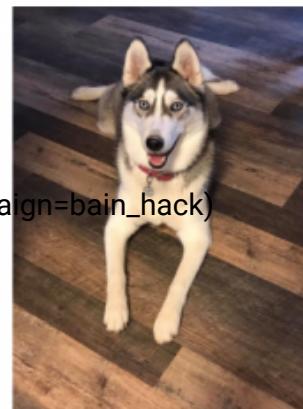
## Women-in-the-loop

A Data Science Hackathon by Bain & Company

28th March–5th April 2020

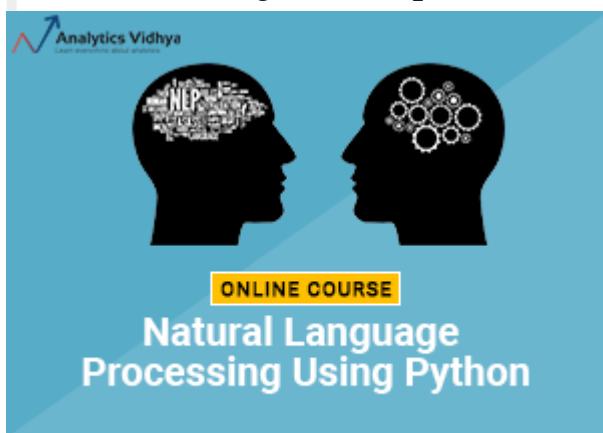
BAIN & COMPANY | Analytics Vidhya

(//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-3.36.47-PM.png>)

*Having an interpretable model, in this case, enables us to test the causality of ty and ultimately can help us to debug the model*



ome with  
nt Starter Kit.



#### Learn to Solve

#### Text Classification Problems Using NLP

Example, we can try to alter our data by adding images of the animals in different backgrounds or simply crop the background out of the images to ensure that the right signal is picked up by our machine learning or language-processing-nlp/?

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>?utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=NLPcourse)



PR's article 12 allows individuals to enquire about algorithmic  
ance industry, questions such as the following can come up and

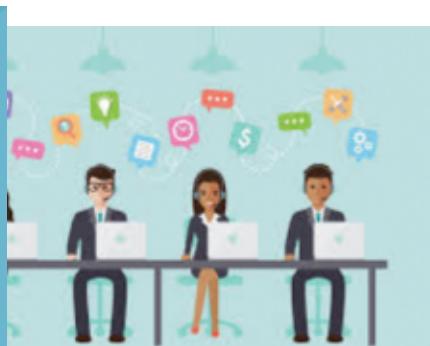
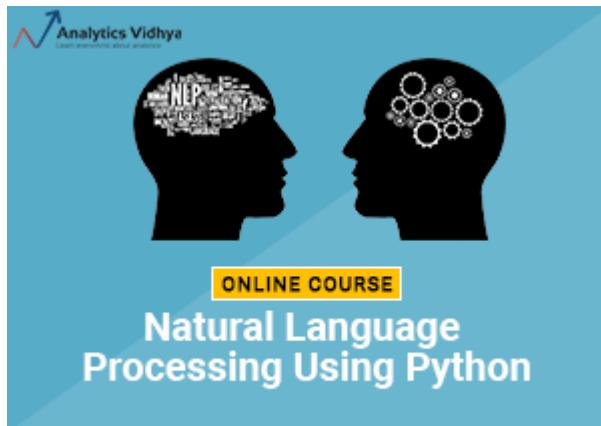
(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>?utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-3.40.35-PM.png>).

## When can we do Away with Interpretability?

Not everything requires interpretability. It is also important to understand when we do *not* need to invest in building interpretable machine learning models.



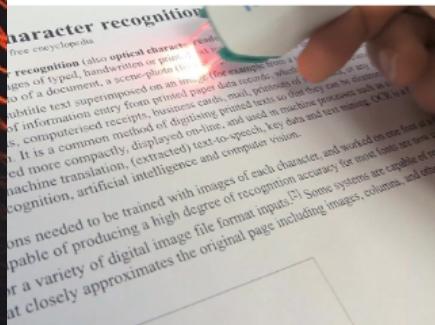
/uploads/2019/08/Screenshot-2019-08-16-at-3.48.28-PM.png)

## Learn to Solve Text Classification Problems Using NLP

Learn to Solve  
Text Classification Problems Using **NLP**

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>) We are good to go as long as we have some internal processes. Say we want to

the end customer. For example, a situation where we are looking to solve natural language processing problems. We want to classify the call recordings since a call resulted in a win or loss as long as we are getting a good performance and it is solving a business problem (e.g., [display&utm\\_campaign=NLPcourse](#))



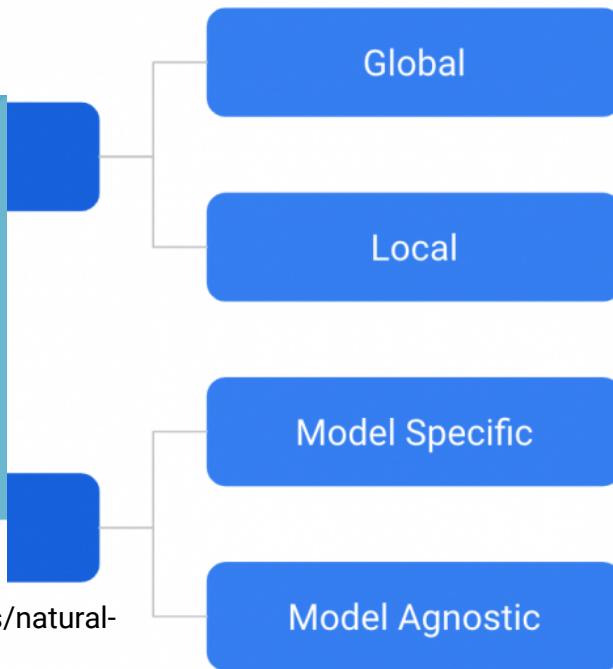
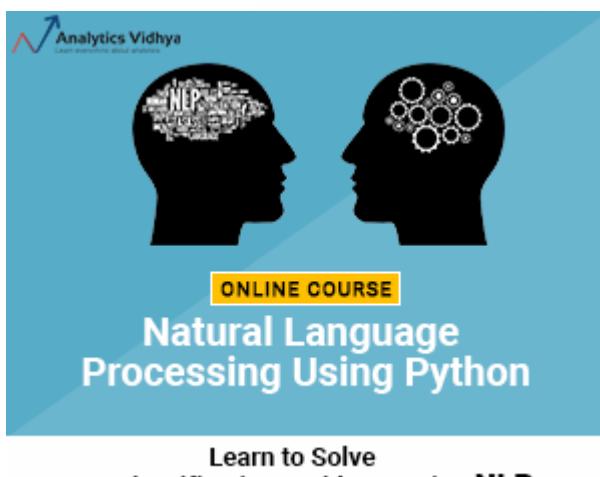
<https://datahack.analyticsvidhya.com/contest/women->

in-the-loop-a-data-science-hackathon-by-  
If the problem is well-studied, we

- If the problem is well studied, we are confident about the results. For example, optical character bain/?  
recognition for which we can get a lot of training data and can rely on a good performance for the task at  
utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)  
hand

# Framework for Interpretable Machine Learning

Now that we have an intuition of what ML interpretability is and why it's important, let's look at the different ways to classify interpretability techniques:



<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>

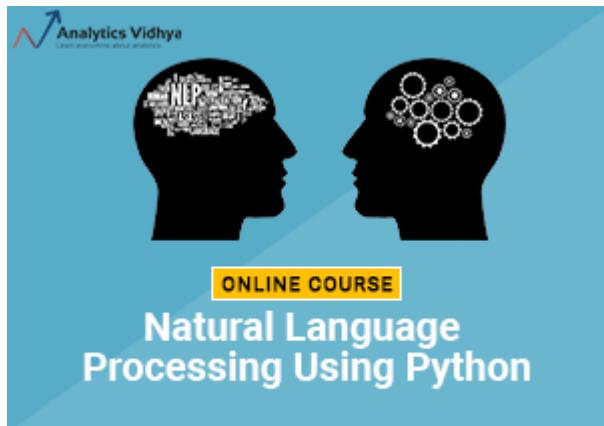
The second way of looking at this is whether we are talking about a technique that works across all types of models (model agnostic) or is tailor made for a particular class of algorithms (model specific).

## Let's Talk About Inherently Interpretable Models

### Linear/Logistic

For linear models such as a linear and logistic regression, we can get the importance from the weights/coefficients of each feature.

Let's revisit that quickly. Suppose we are trying to predict an employee's salary using linear regression. The independent variables are experience in years and a previous rating out of 5.



**Learn to Solve  
Text Classification Problems Using NLP**

([http://courses.analyticsvidhya.com/courses/natural-language-processing.html?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=NLP\\_course](http://courses.analyticsvidhya.com/courses/natural-language-processing.html?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=NLP_course))

## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

BAIN & COMPANY |

(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

$$1 * \text{experience} + W2 * \text{rating}$$

([https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=bain\\_hack](https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack))

Finally tell us whether the experience is more important towards salary or rating. This is a very basic technique that can be used for both global and local

on in the below articles:

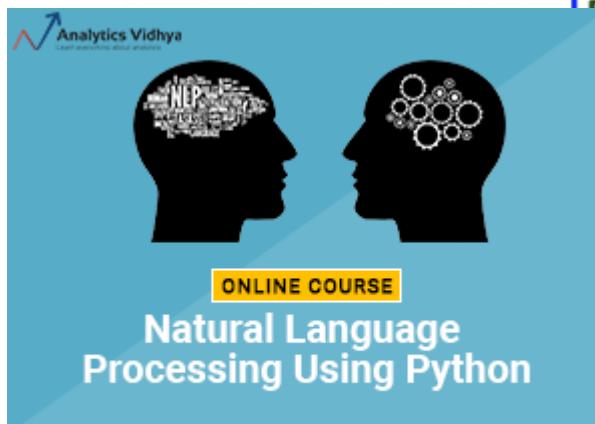
([https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=NEP\\_course](https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=NEP_course))

([https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=NEP\\_course](https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=NEP_course))

R and Python

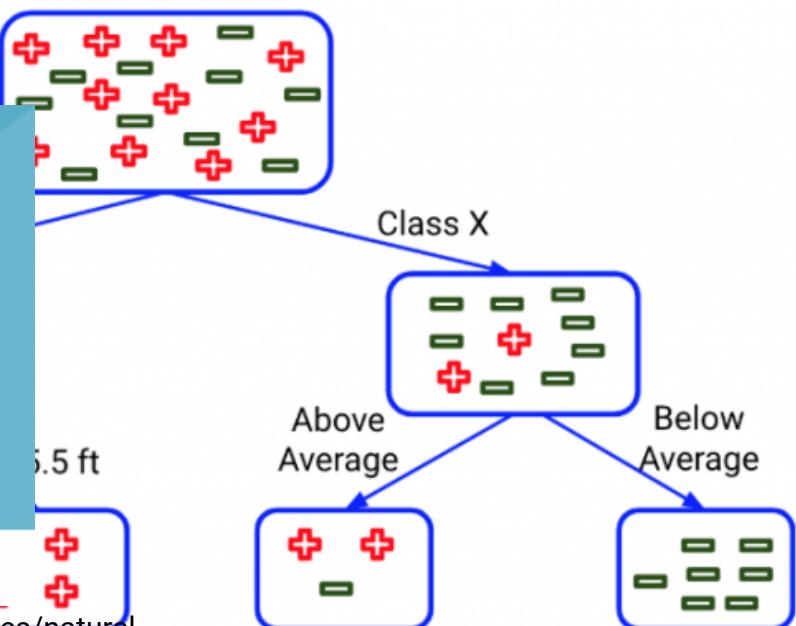
([https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=NEP\\_course](https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=NEP_course))

([https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=NEP\\_course](https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=NEP_course))

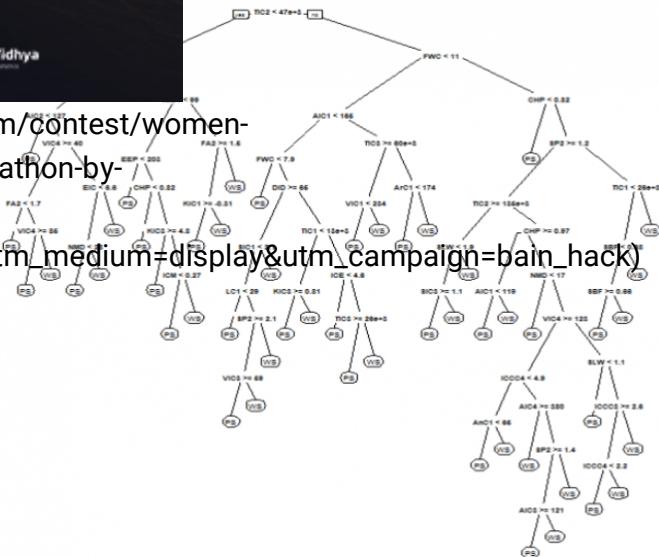


Learn to Solve  
Text Classification Problems Using NLP

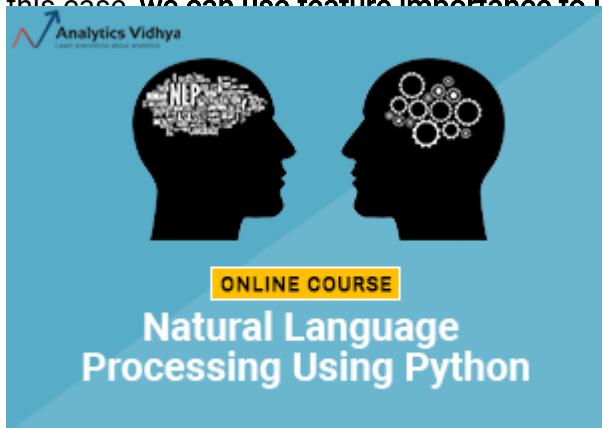
(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)  
utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)



For a small decision tree, we can use the above diagram. However, if we have a lot of features and we are training a deep decision tree for, let's say, a depth of 8 or 9, there will be too many decision rules to present effectively. In this case, we can use feature importance to interpret the importance of each feature at a global level.



### Learn to Solve Text Classification Problems Using NLP

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

- Step 1: Go through all the splits in which the feature was used

to maximize the decrease in impurity (Gini/information gain) compared to the parent node weighted by its importance as they help us incorporate the number of samples well into the equation. The formula is:

$$\frac{N_{Right}}{N_{parent}} \cdot Gini_{Right} - \frac{N_{Left}}{N_{parent}} \cdot Gini_{Left}$$



(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

- Step 2: Calculate the weighted sum of Gini impurity for each feature

$\sum \left( \frac{N_{Right}}{N_{parent}} \cdot Gini_{Right} - \frac{N_{Left}}{N_{parent}} \cdot Gini_{Left} \right)$

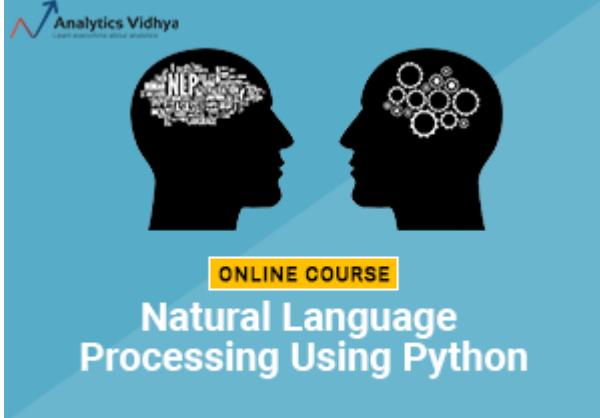
(<https://www.analyticsvidhya.com/uploads/2019/08/Screenshot-2019-08-16-at-4.12.05-PM.png>)

Here, again, this is a model-specific technique that can be used for only global explanations. This is because we are looking at the overall importance and not at each prediction.

Learn more about decision trees in [this superb tutorial](#)

([https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/?utm\\_source=blog&utm\\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python](https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/?utm_source=blog&utm_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python)).

## An Example of Feature Importance



(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



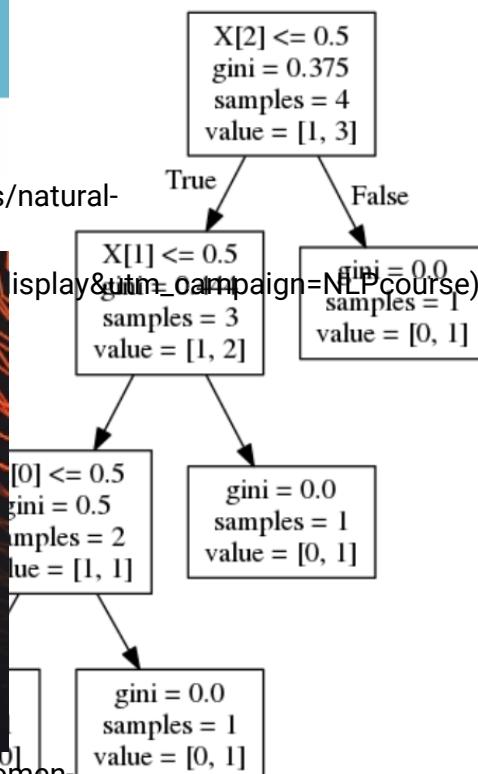
(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

Since each feature is used once in our case, there is no need to calculate the sum.

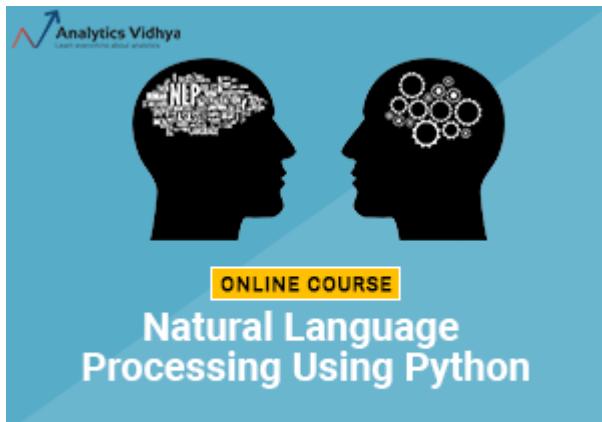
- For X[2] :
  - feature\_importance =  $(4 / 4) * (0.375 - (0.75 * 0.444)) = 0.042$
- For X[1] :
  - feature\_importance =  $(3 / 4) * (0.444 - (2/3 * 0.5)) = 0.083$
- For X[0] :
  - feature\_importance =  $(2 / 4) * (0.5) = 0.25$

ample. This will help you visualize what we've covered so far (and

roles and Gini impurity values as shown in the below figure. Each directly use the formula to calculate feature importance:



Take a moment to pause here and calculate this on your own. You will have a much better grasp on the concept once you do it by yourself.



Learn to Solve  
Text Classification Problems Using **NLP**

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/2>) *Future Imp.* - *Sum of*

## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

(//datacamp.com/ensemble-learning-with-python-codes)

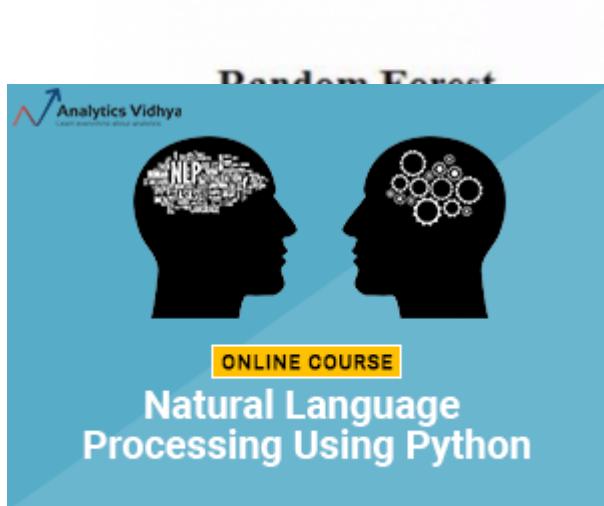
in-the- (<https://www.analyticsvidhya.com/blog/2018/06/comprehensive-guide-for-ensemble-models/>)

bain/?utm\_source=blog&utm\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-  
utm\_medium=python)anner2&utm\_medium=display&utm\_campaign=bain\_hack)

## Model Agnostic Techniques

So far, we have discussed model-specific techniques for linear and logistic regression, as well as decision trees. We also spoke about the feature importance methods that are used for ensemble methods. I'm sure you're wondering – what about other models?

Now, we know that some models are hard to interpret, such as random forest and gradient boosting.



Learn to Solve  
Text Classification Problems Using **NLP**

([http://courses.analyticsvidhya.com/courses/natural-](http://courses.analyticsvidhya.com/courses/natural-language-processing/)

(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-6.14.59-PM.png>). We

Yes. But it does not tell us whether a particular feature affects the display&utm\_medium=NLP course) RY important in certain cases.

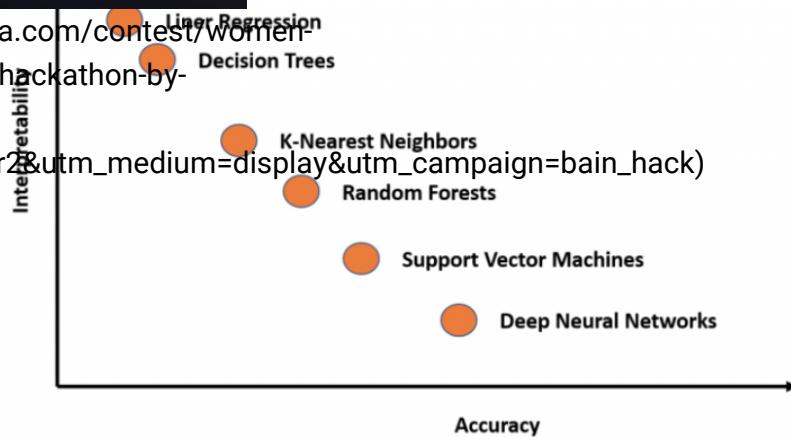
## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

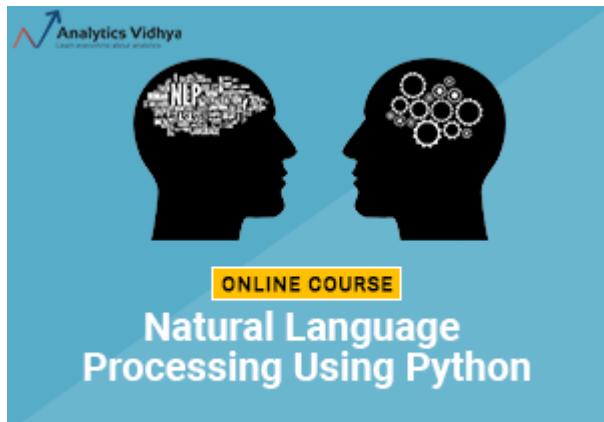
BAIN & COMPANY |

(<http://datahack.analyticsvidhya.com/contests/women-in-the-loop-a-data-science-hackathon-by-bain/>)?



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-6.20.59-PM.png>).

I really like this plot. As the complexity of the machine learning model increases, we get better performance but lose out on interpretability. You should keep this figure handy the next time you're building your own model.

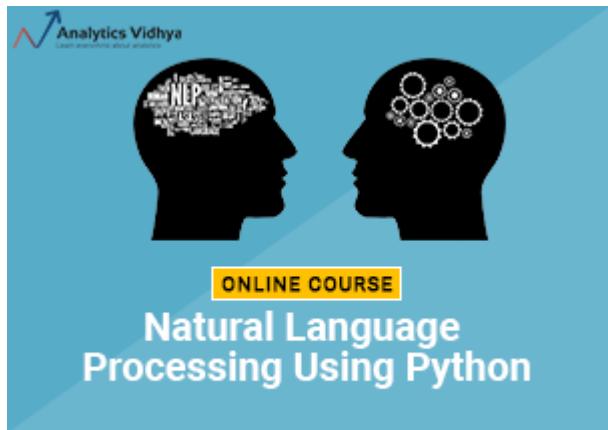


[\(http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)

Interpretability is just another aspect of model accuracy. We capture the world by collecting data, and then we can use that data to predict the world with a machine learning model. Interpretability is just another way to look at a black box using a simpler, more interpretable model.



[\(/datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm\\_source=Sticky\\_banner2&utm\\_medium=display&utm\\_campaign=bain\\_hack\)](http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)



Learn to Solve  
Text Classification Problems Using **NLP**

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

## Global Surrogate Method

The first model agnostic method we will discuss here is the global surrogate method. **A global surrogate model is an interpretable model that is trained to approximate the predictions of a black-box model.**

We can draw conclusions about the black-box model by interpreting the surrogate model. So, we are basically solving machine learning interpretability by using more machine learning!

For example, we could interpret a random forest classifier using a simple decision tree to explain its predictions:

**Complex Model**




**ONLINE COURSE**

## Natural Language Processing Using Python

Learn to Solve  
Text Classification Problems Using **NLP**

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)  
(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-6.28.49-PM.png>).

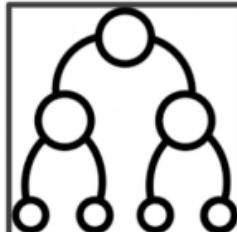
**Women-in-the-loop**

A Data Science Hackathon by Bain & Company

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

**Simpler Model**



Decision Tree Classifier  
Predictions: [0,0,0,1,1,0]

racy: 83.33 % accuracy

isplay&utm\_campaign=NLPcourse)  
predictions of the black-box model (which is a random forest in our  
uracy, we can use it to explain the random forest classifier.

*interpretable surrogate model that is trained to  
s of a black-box model and draw conclusions.*

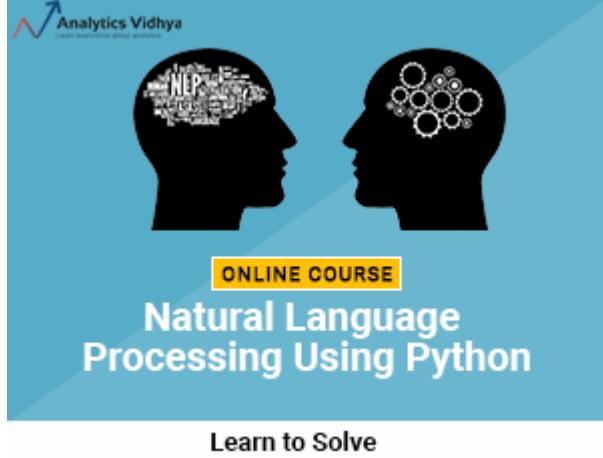
nd how a global surrogate model works:

- (<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)
1. We get predictions from the black-box model
  2. Next, we select an interpretable model (Linear, decision tree, etc.)
  3. We train an interpretable model on the original dataset and use black box predictions as the target
  4. Measure the performance of the surrogate model
  5. Finally, we interpret the surrogate model to understand how the black-box model is making its decisions

## LIME (Local Interpretable Model agnostic Explanations)

The global surrogate method is good for looking at an interpretable model that can explain predictions for a black-box approach. However, this will not work well if we want to understand how a single prediction was made for a given observation.

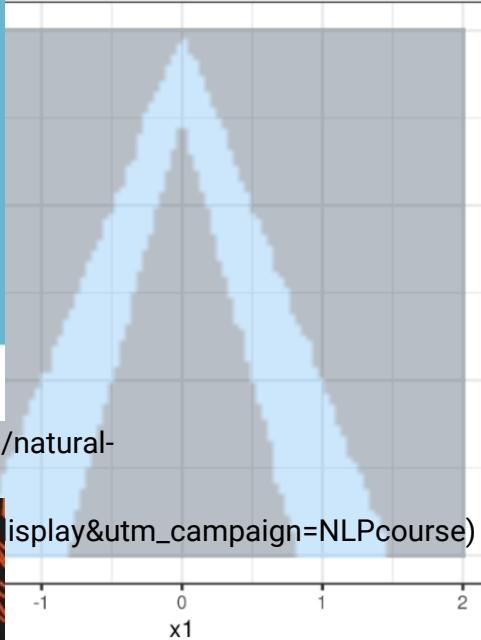
This is where we use the LIME technique which stands for local interpretable model agnostic explanations. LIME is based on the work presented in [this paper](https://arxiv.org/abs/1602.04938) (<https://arxiv.org/abs/1602.04938>). Let us understand how LIME works using an example.



(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)



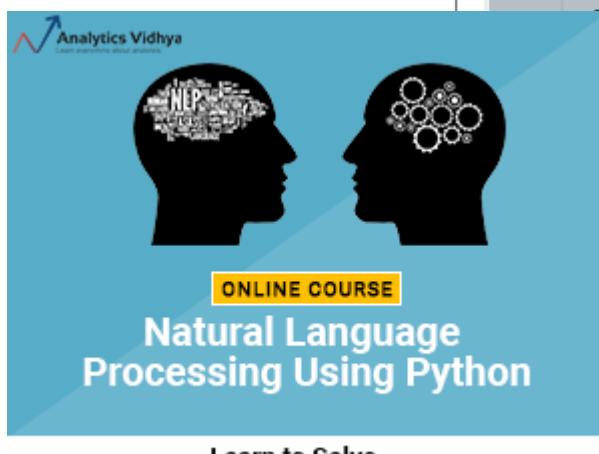
(<http://uploads/2019/08/Screenshot-2019-08-16-at-6.45.28-PM.png>).

ation problem. As we can see in the above image, we have a two features.

values of  $x_1$  and  $x_2$  for the observation in yellow (in the below image).tribution to generate fake data around the observation:

(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)





(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

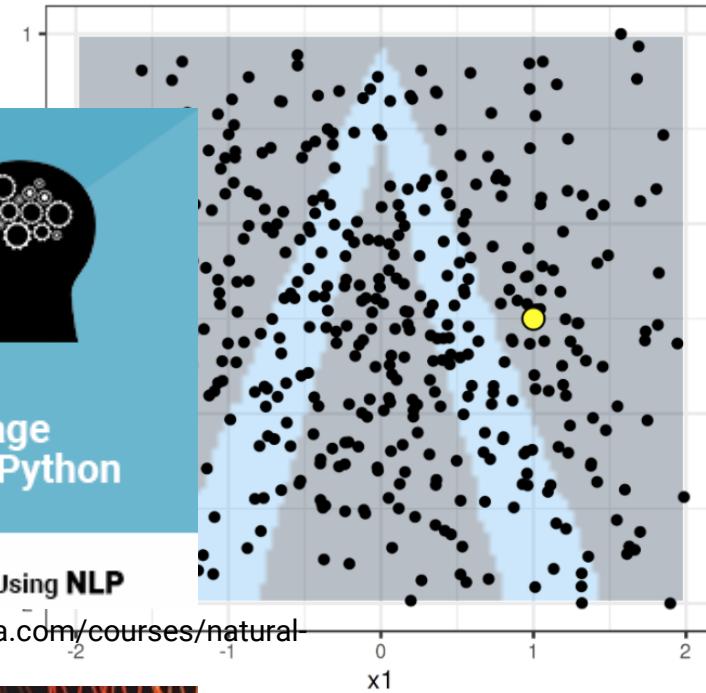
## Women-in-the-loop

A Data Science Hackathon by Bain & Company

28th March–5th April 2020

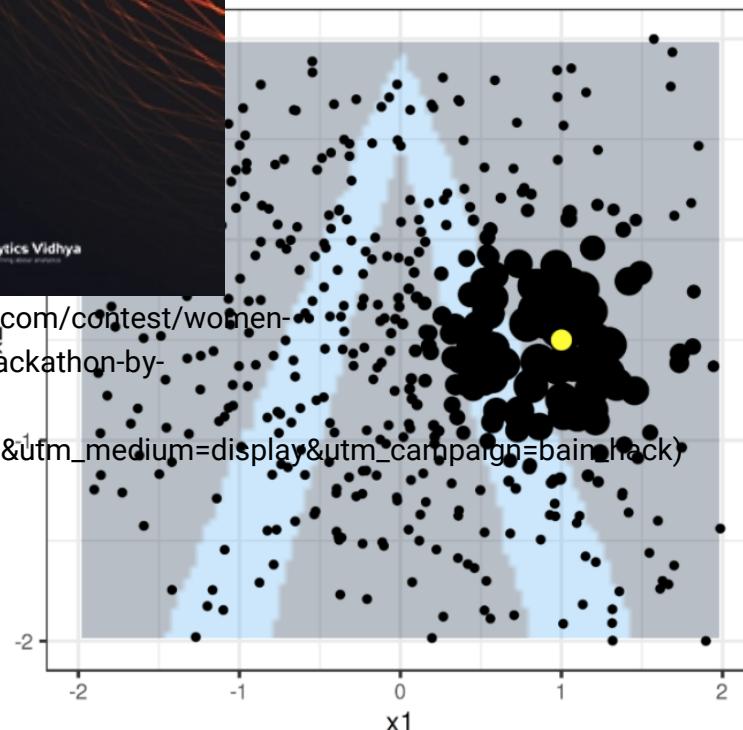
BAIN & COMPANY |

(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)  
 utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-6.48.17-PM.png>).

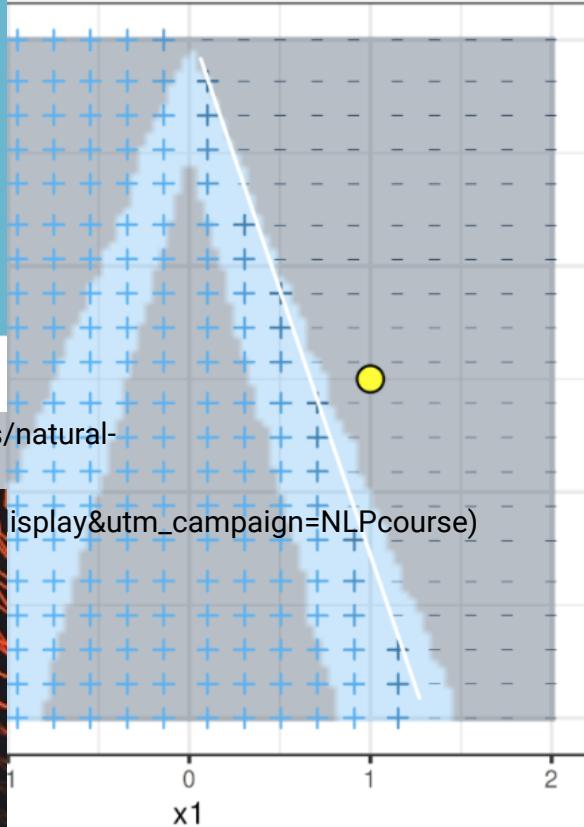
that are closer to our observation:



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-6.50.42-PM.png>).

We train an interpretable model over the fake data generated from the distribution. Now, we have a new local decision boundary for the locally learned model (in white) that can be used to understand the contributions of  $x_1$  and  $x_2$  towards the prediction of our observation:

The image contains two banners. The top banner is for an 'ONLINE COURSE' titled 'Natural Language Processing Using Python' offered by Analytics Vidhya. It features two stylized human head profiles facing each other, with one head containing gears and the other containing text. Below the title is the text 'Learn to Solve Text Classification Problems Using NLP'. The URL [\(http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/) is also present. The bottom banner is for a 'Women-in-the-loop' Data Science Hackathon by Bain & Company. It features a dark background with red abstract lines and the text 'Women-in-the-loop', 'A Data Science Hackathon by Bain & Company', and '28th March–5th April 2020'. The URL [\(http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/\)](http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/) is also present.



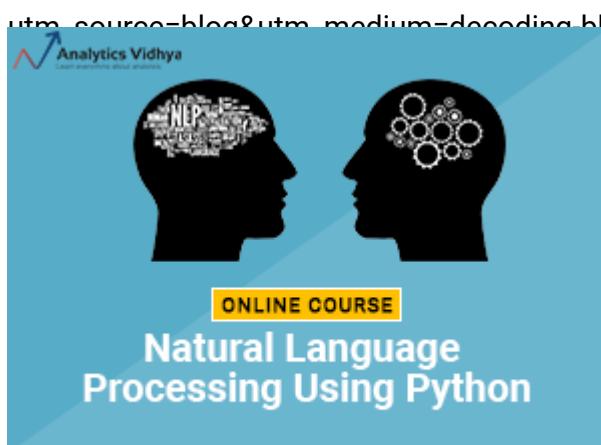
[https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm\\_source...](https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm_source...)

## Python Implementation of Interpretable Machine Learning Techniques

My favorite part of the article – building interpretable machine learning models in Python!



Here, we will work on the implementation of both the methods we covered above. We will use the [big mart sales problem \(\[https://datahack.analyticsvidhya.com/contest/practice-problem-big-mart-sales-iii/?utm\\\_source=blog&utm\\\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python\]\(https://datahack.analyticsvidhya.com/contest/practice-problem-big-mart-sales-iii/?utm\_source=blog&utm\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python\)\)](https://datahack.analyticsvidhya.com/contest/practice-problem-big-mart-sales-iii/?utm_source=blog&utm_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python)



(<http://courses.analyticsvidhya.com/courses/natural-language-processing-and-understanding-interpretable-machine-learning-models>)

## Building and Understanding Interpretable Machine Learning Models

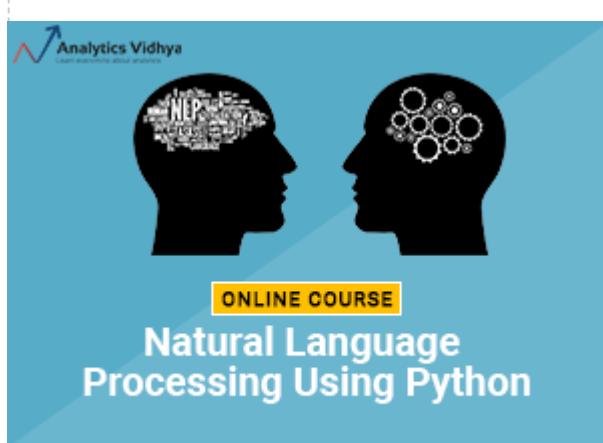
([http://courses.analyticsvidhya.com/courses/big-mart-sales-prediction?utm\\_source=blog&utm\\_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python](http://courses.analyticsvidhya.com/courses/big-mart-sales-prediction?utm_source=blog&utm_medium=decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python)) to fully understand how to build models using this interpretability part.



```
//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
from xgboost.sklearn import XGBRegressor
from sklearn.preprocessing import OneHotEncoder, LabelEncoder
from sklearn import tree

import matplotlib.pyplot as plt
%matplotlib inline
```

## Reading Data



```
by median and Outlet_Size with mode
].median(), inplace=True)
].mode()[0], inplace=True)
```

**Learn to Solve**  
**Text Classification Problems Using NLP**  
**feature engineering**  
[\(http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)



```
display&utm_campaign=NLPcourse)
classifier'].apply(lambda df: df[0:2])
_Combined'].map({'FD':'Food', 'NC':'Non-Consumable', 'DR':'Dr
```

<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain-and-company-modifying-categories-of-item-fat-content>

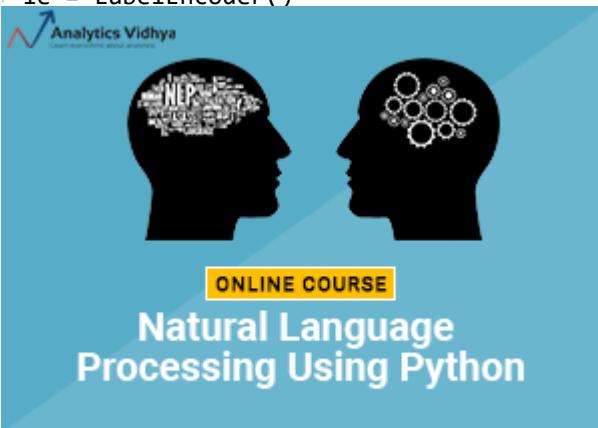
```
bain['Item_Fat_Content'] = df['Item_Fat_Content'].replace({'LF':'Low Fat', 'reg':'Regular', 'low fat':'Low Fat'})
utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)

df['Item_Fat_Content'].value_counts()
```

## Data Preprocessing

```
# label encoding the ordinal variables
```

```
le = LabelEncoder()
```



Learn to Solve  
Text Classification Problems Using NLP

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

```
'Item_Type_Combined', 'Outlet'])
```

```
display&utm_campaign=NLPcourse)
```

## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

BAIN & COMPANY | 

(/ydatahack@analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?)

Creating the training and validation set

```
utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.25, random_state=42)
```

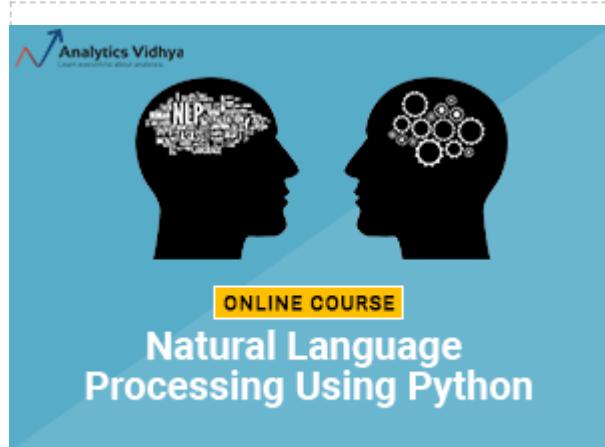
## Training a Decision Tree Model

```
dt = DecisionTreeRegressor(max_depth = 5, random_state=10)
```

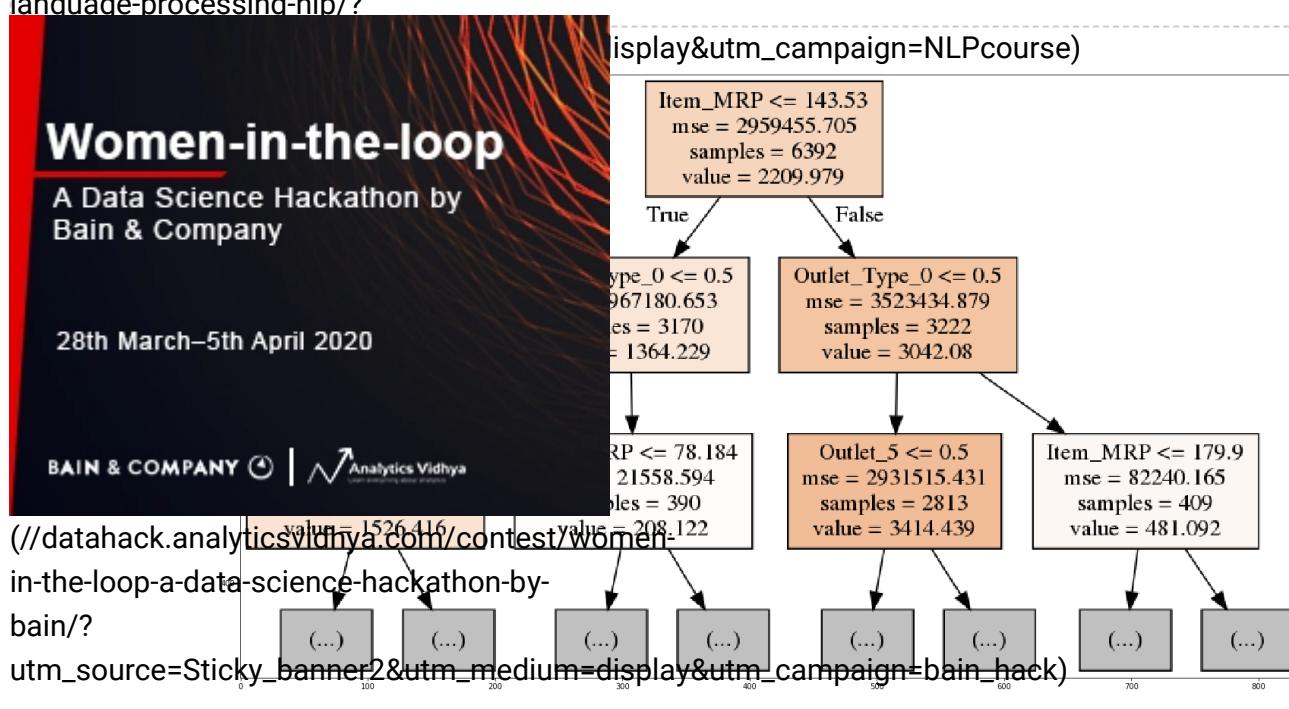
```
# fitting the decision tree model on the training set
```

```
dt.fit(X_train, y_train)
```

## Use the Graphviz library to visualize the decision tree



```
out_file='tree.dot', feature_names=X_train.columns, filled=T
```

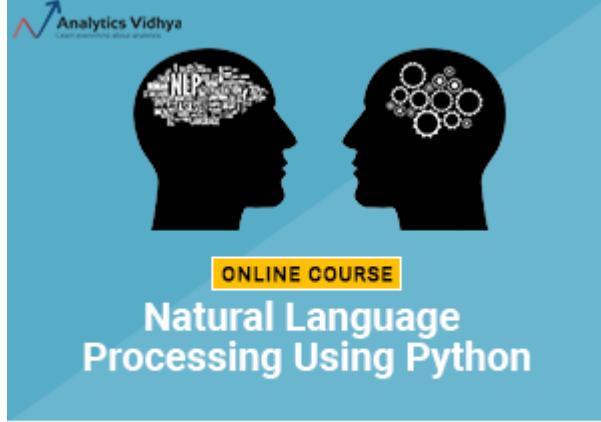


(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/download.png>).

This visualization of our decision tree clearly displays the rules it is using to make a prediction. Here, **Item\_MRP & Outlet\_Type** are the first features that are affecting the sales of various items at each outlet. If you want to look at the complete decision tree, you can easily do that by changing the *max\_depth* parameter using the *export\_graphviz* function.

## Feature Importance

Now, we will have a look at the feature importance for each feature in case of a random forest.



Learn to Solve  
Text Classification Problems Using **NLP**

([http://courses.analyticsvidhya.com/courses/natural-](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)

[language-processing-nlp/?values\\_by='importance', ascending=False\).head\(10\) display&utm\\_campaign=NLPcourse](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?values_by='importance', ascending=False).head(10) display&utm_campaign=NLPcourse))



(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

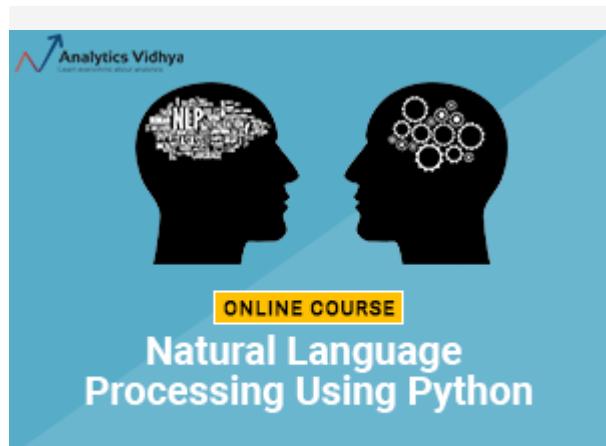
variable	importance
Item_MRP	0.552067
Outlet_Type_0	0.322084
Outlet_Type_3	0.046778
Outlet_Years	0.043134
Outlet_5	0.034728
Item_Visibility	0.000638
Item_Weight	0.000279
Outlet_Type_1	0.000113
Outlet_Location_Type_2	0.000082
Outlet_Size_1	0.000060

(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/Screenshot-2019-08-16-at-7.24.15-PM.png>).

The random forest model gives a similar interpretation. Item\_MRP still remains the most important feature (exactly as the decision tree model above). **Relative importance also helps us compare each feature.** For example, Outlet\_Type\_0 is a much more important feature than other outlet types.

## Exercise

As an exercise try calculating the feature importance for the decision tree we fit earlier and compare:



**ONLINE COURSE**

## Natural Language Processing Using Python

Learn to Solve Text Classification Problems Using NLP

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

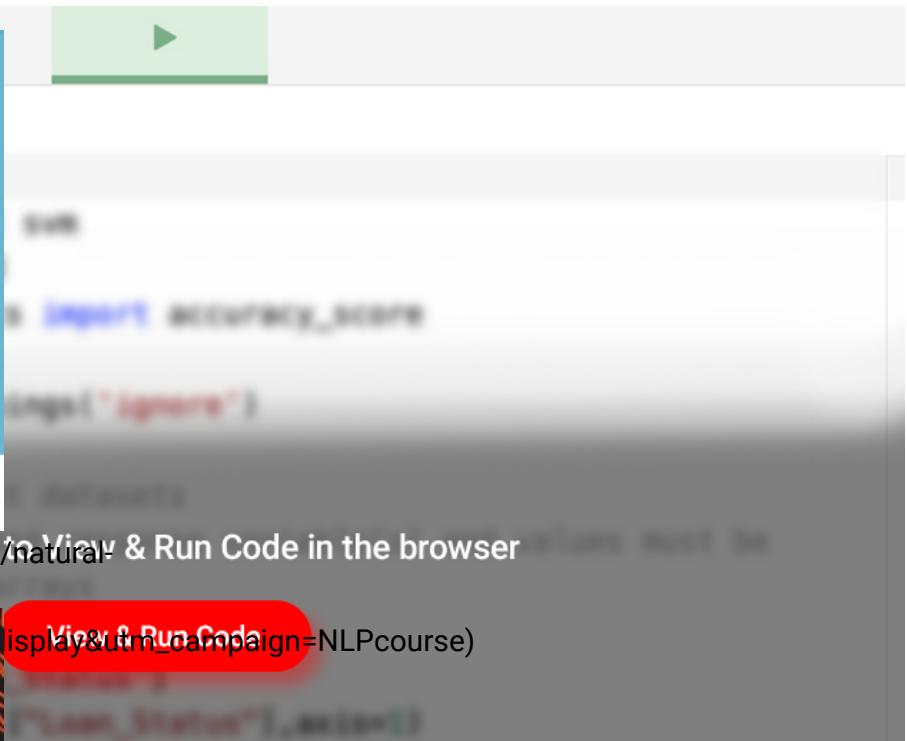
[View & Run Code](#)

**Women-in-the-loop**

A Data Science Hackathon by Bain & Company

28th March–5th April 2020

BAIN & COMPANY | 



Next, we will create a surrogate decision tree model for this random forest model and see what we get.

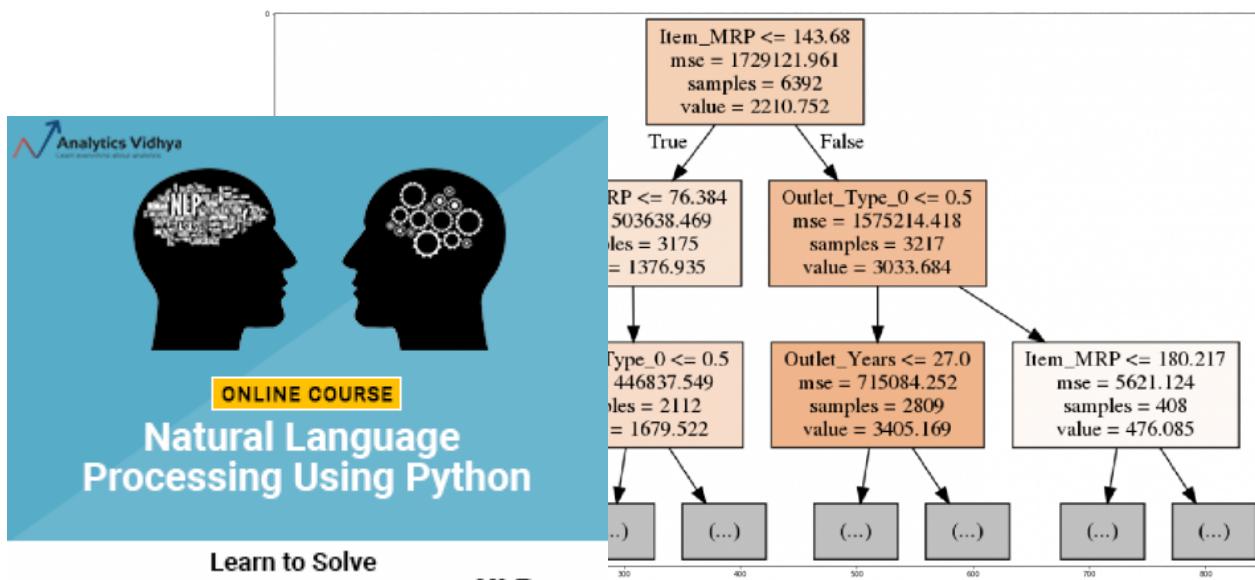
```
ext=https://www.analyticsvidhya.com/blog/2019/08/decoding-machine-learning-models-python/?utm_source=blog&utm_medium=6-models?utm_source=coding-window-blog&source=coding-window-
```

<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>

```
# saving the predictions of Random Forest as new target
new_target = rf.predict(X_train)
```

```
# defining the interpretable decision tree model
dt_model = DecisionTreeRegressor(max_depth=5, random_state=10)

# fitting the surrogate decision tree model using the training set and new target
dt_model.fit(X_train,new_target)
```



Learn to Solve Text Classification Problems Using NLP

(<https://cdn.analyticsvidhya.com/wp-content/uploads/2019/08/a.png>).

[language-processing-nlp/?](https://language-processing-nlp/)



(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)  
utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

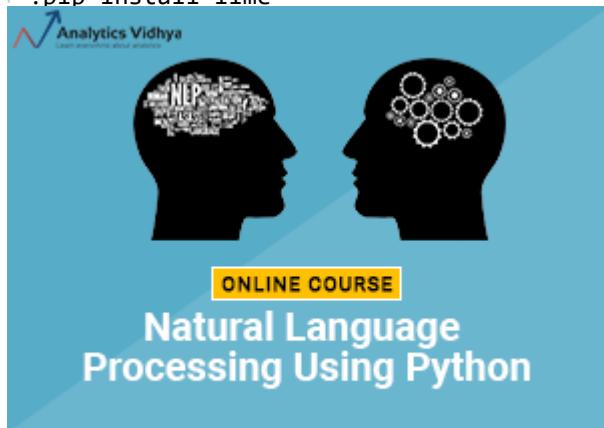
target and can be used as a surrogate model to explain the display&utm\_campaign=NEPcourse). Finally, we can use it for any other complex model. Just make sure you get wrong interpretations (a nightmare!).

## Generate local interpretations of black-box models

In R & Python using the [LIME package](#)  
help into implementation for the same to check the local E:

```
# installing lime library
```

```
!pip install lime
```



### Learn to Solve Text Classification Problems Using NLP

```
(#creating the explainer function
explainer = lime.LimeExplainer(X_train.values, mode="regression", feature_names=X_train.columns)
display&utm_campaign=NLPcourse)
```

## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

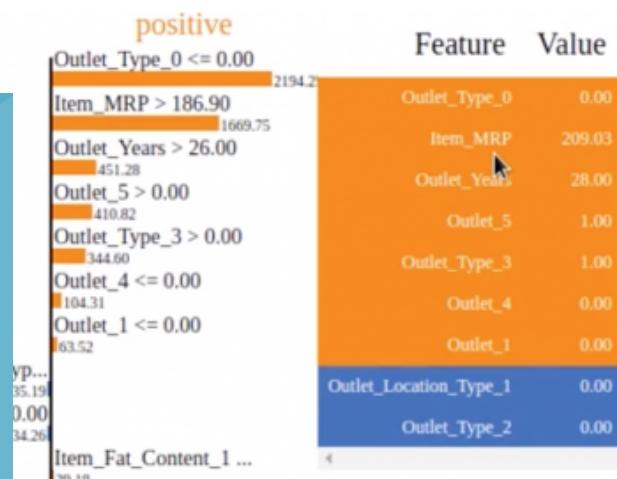
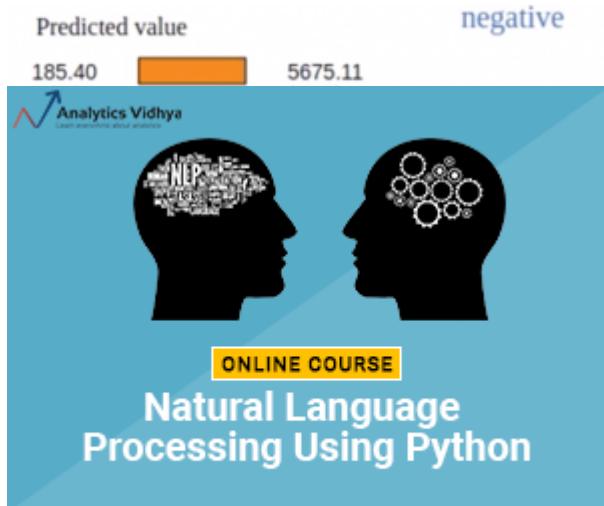
BAIN & COMPANY | Analytics Vidhya

```
observation)[0]}

0328626019
/uploads/2019/08/Screenshot-2019-08-16-at-8.05.48-PM.png).
(//datahack.analyticsvidhya.com/contest/women-
in-the-loop-a-data-science-hackathon-by-
bain/?
utm_source=Sticky_banner2&utm_medium=display&utm_campaign=bain_hack)
```

## Generate Explanations using LIME

```
# explanation using the random forest model
explanation = explainer.explain_instance(X_observation.values[0], rf_model.predict)
explanation.show_in_notebook(show_table=True, show_all=False)
print(explanation.score)
```



Learn to Solve  
Text Classification Problems Using **NLP**  
blue signifies the positive impact and  
example, *Item\_MRP* has a positive impact on sales.  
language-processing-nlp/?



(//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

The image shows a screenshot of a Jupyter notebook interface from DataCamp. The code cell contains the following Python code:

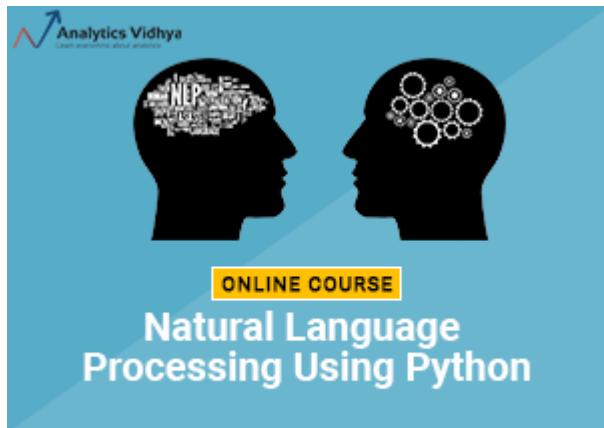
```

trainhead.read_csv('train.csv')
train_y = train['Label'].values
train_x = train.drop(['Label'], axis=1)
train_x = train_x.fillna(0)
train_x = train_x.values
train_x[0]

```

Below the code cell, there's a red button labeled "View & Run Code". Above the code cell, there's a message "Login/ Signup to View & Run Code in the browser".

([https://id.analyticsvidhya.com/auth/login/?next=https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm\\_source=blog&utm\\_medium=6-](https://id.analyticsvidhya.com/auth/login/?next=https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?utm_source=blog&utm_medium=6-)



Learn to Solve  
Text Classification Problems Using **NLP**  
(<https://dileepreddy.analyticsvidhya.com/courses/natural-language-processing-nlp/>)



(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>).  
utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

Share this:

[\(https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=linkedin&nb=1\)](https://www.linkedin.com/shareArticle?utm_source=blog_article&utm_medium=blog&utm_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2315-1)

[\(https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=facebook&nb=1\)](https://www.facebook.com/sharer/sharer.php?u=https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=facebook&nb=1)

[\(https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=twitter&nb=1\)](https://twitter.com/intent/tweet?url=https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=twitter&nb=1)

(<https://www.analyticsvidhya.com/blog/2019/08/decoding-black-box-step-by-step-guide-interpretable-machine-learning-models-python/?share=pocket&nb=1>)



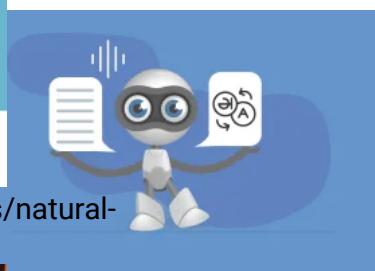
## ONLINE COURSE

# Natural Language Processing Using Python

Learn to Solve  
Text Classification Problems Using **NLP**

[\(http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)

oding-black-box-step-by-step-guide-interpretable-machine-learning-models-



# Women-in-the-loop

A Data Science Hackathon by Bain & Company

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

isplay&utm\_campaign=NLPcourse)  
/2019/03/06/analyticsvidhya.com/blog/2019/03/06/analyticsvidhya.com/blog/  
ral-language-processing-nlp-  
essions-datahack-summit-  
engaluru/).  
ing NLP Hack Sessions to  
out for at DataHack Summit  
<http://www.analyticsvidhya.com/blog/2019/10/top-7-natural-language-processing-nlp-hack-sessions-datahack-summit-2019-bengaluru/>

(<https://www.analyticsvidhya.com/blog/2019/10/top-7-natural-language-processing-nlp-hack-sessions-datahack-summit-2019-bengaluru/>)

[data-science-projects-github-showcase-your-skills/](https://www.analyticsvidhya.com/blog/2019/10/7-data-science-projects-github-showcase-your-skills/)).

Here are 7 Data Science Projects on GitHub to Showcase your Machine Learning Skills!

(<https://www.analyticsvidhya.com/blog/2019/09/7-data-science-projects-github-showcase-your-skills/>)

September 2, 2019

In "Advanced"

(<https://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

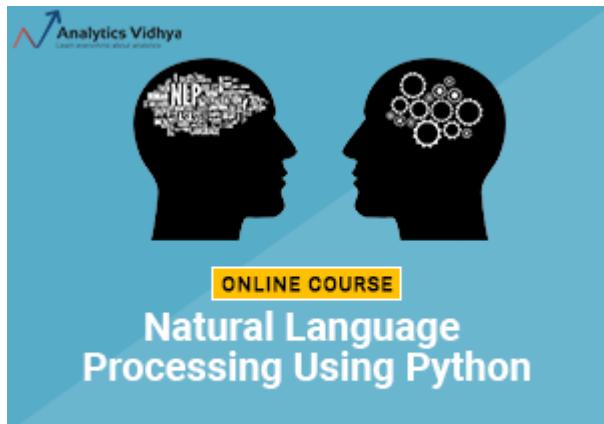
utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

Save money with the  
Smart Light Starter Kit.



**TAGS :** [INTERPRETABLE MACHINE LEARNING](https://www.analyticsvidhya.com/blog/tag/interpretable-machine-learning/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/INTERPRETABLE-MACHINE-LEARNING/](https://www.analyticsvidhya.com/blog/tag/interpretable-machine-learning/)), [INTERPRETABLE ML](https://www.analyticsvidhya.com/blog/tag/interpretable-ml/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/INTERPRETABLE-ML/](https://www.analyticsvidhya.com/blog/tag/interpretable-ml/)), [LIME](https://www.analyticsvidhya.com/blog/tag/lime/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/LIME/](https://www.analyticsvidhya.com/blog/tag/lime/)), [LIVE CODING](https://www.analyticsvidhya.com/blog/tag/live-coding/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/LIVE-CODING/](https://www.analyticsvidhya.com/blog/tag/live-coding/)), [MACHINE LEARNING](https://www.analyticsvidhya.com/blog/tag/machine-learning/) ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MACHINE-LEARNING/](https://www.analyticsvidhya.com/blog/tag/machine-learning/)), [MACHINE LEARNING](https://www.analyticsvidhya.com/blog/tag/machine-learning/)

[INTERPRETABILITY \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MACHINE-LEARNING-INTERPRETABILITY/\)](https://www.analyticsvidhya.com/blog/tag/machine-learning-interpretability/), [PYTHON \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/PYTHON/\)](https://www.analyticsvidhya.com/blog/tag/python/)



Learn to Solve  
Text Classification Problems Using **NLP**

(<https://www.analyticsvidhya.com/blog/2019/08/11-data-visualizations-python-r-tableau-d3js/>)

display&utm\_campaign=NLPcourse)

## Women-in-the-loop

A Data Science Hackathon by  
Bain & Company

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

(//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/)

Ankit Choudhary (<https://www.analyticsvidhya.com/blog/author/Ankit2106/>)  
IIT Bombay Graduate with a Masters and Bachelors in Electrical Engineering. I have previously worked as a lead decision scientist for Indian National Congress deploying statistical models (Segmentation, K-Nearest Neighbours) to help party leadership/Team make data-driven decisions. My interest lies in putting data in heart of business for data-driven decision making.

in\_ (<https://www.linkedin.com/in/ankit-choudhary-b9360826/>)

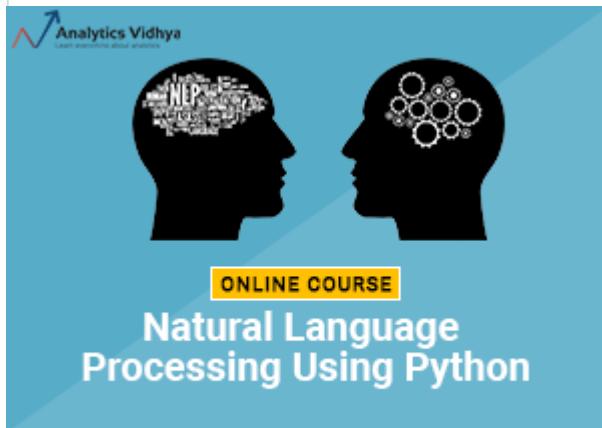


(<https://www.analyticsvidhya.com/blog/author/Ankit2106/>)

## LEAVE A REPLY

Your email address will not be published.

## Comment



Learn to Solve  
Text Classification Problems Using **NLP**  
Email required  
[\(http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/?\)](http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/)

display&utm\_campaign=NLPcourse)

An advertisement for the "Women-in-the-loop" Data Science Hackathon. It features a dark background with red abstract lines. The text includes: "Women-in-the-loop", "A Data Science Hackathon by Bain & Company", "28th March–5th April 2020", "BAIN & COMPANY | Analytics Vidhya", "Instantly Convert EDW in-the-loop-a-data-science-hackathon-by-bain.com", "90% Automatic Conversion", and "4x Speed vs. Traditional Appro".

BAIN & COMPANY | Analytics Vidhya

//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/

Instantly Convert EDW in-the-loop-a-data-science-hackathon-by-bain.com

90% Automatic Conversion

4x Speed vs. Traditional Approach

CONTACT US

IMPETUS

## POPULAR POSTS




**ONLINE COURSE**

### Natural Language Processing Using Python

Learn to Solve Text Classification Problems Using **NLP**

(<https://www.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

Naïve Bayes Algorithm with codes in Python and R

(<https://www.analyticsvidhya.com/courses/naive-bayes-explained/>)

for (aspiring) data scientists

(<https://www.analyticsvidhya.com/courses/introductory-guide-on-linear-programming-explained-in-python/>)

XGBoost with codes in Python

(<https://www.analyticsvidhya.com/courses/complete-guide-parameter-tuning-xgboost-with-codes-in-python/>)

BAIN & COMPANY | Analytics Vidhya

(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop/>)

**6 Python Libraries to Interpret Machine Learning Models and Build Trust**

(<https://www.analyticsvidhya.com/blog/2020/03/6-python-libraries-interpret-machine-learning-models/>)

MARCH 25, 2020 Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

**Using Graphs to Identify Social Media Influencers** (<https://www.analyticsvidhya.com/blog/2020/03/using-graphs-to-identify-social-media-influencers/>)

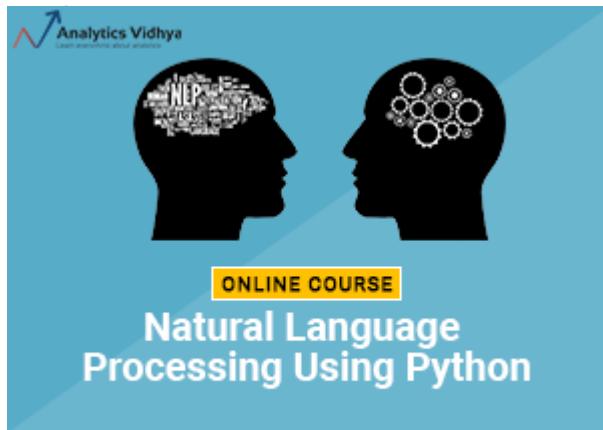
MARCH 24, 2020

**Coronavirus Analysis: Will Social Distancing Help Prevent the Spread?**

(<https://www.analyticsvidhya.com/blog/2020/03/coronavirus-analysis-will-social-distancing-help-prevent-the-spread/>)

MARCH 23, 2020

**Free GPUs for Everyone! Get Started with Google Colab for Machine Learning and Deep Learning**  
(<https://www.analyticsvidhya.com/blog/2020/03/google-colab-machine-learning-deep-learning/>)



Learn to Solve  
Text Classification Problems Using **NLP**

(<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

(utm\_medium=display&utm\_campaign=NLPcourse)



(<http://datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/>)

utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)



**Analytics Vidhya** (<http://courses.analyticsvidhya.com/courses/natural-language-processing-nlp/>)

About Us (<https://www.analyticsvidhya.com/about-me/>)

Our Team (<https://www.analyticsvidhya.com/about-me/team/>)

## Women-in-the-loop

Careers (<https://www.analyticsvidhya.com/about-me/career-analytics-Discussions>) (<https://discuss.analyticsvidhya.com/>)

### Data Science Hackathon by Bain & Company

Contact us (<https://www.analyticsvidhya.com/contact/>)

28th March–5th April 2020

BAIN & COMPANY | Analytics Vidhya

(//datahack.analyticsvidhya.com/contest/women-in-the-loop-a-data-science-hackathon-by-bain/?utm\_source=Sticky\_banner2&utm\_medium=display&utm\_campaign=bain\_hack)

f (<https://www.linkedin.com/company/analytics-vidhya/>)

(https://www.linkedin.com/company/analytics-vidhya/)

© Copyright 2013-2020 Analytics Vidhya

[Privacy Policy](#) [Terms of Use](#) [Refund Policy](#)



(https://apps.apple.com/us/app/analytics-vidhya/id1470025572)

## Data Science

Blog (<https://www.analyticsvidhya.com/blog/>)

Hackathon (<https://datahack.analyticsvidhya.com/>)

Apply Jobs (<https://www.analyticsvidhya.com/jobs/>)

## Companies

Post Jobs (<https://www.analyticsvidhya.com/corporate/>)

Trainings (<https://courses.analyticsvidhya.com/>)

Hiring Hackathons (<https://datahack.analyticsvidhya.com/>)

Advertising (<https://www.analyticsvidhya.com/contact/>)

x

-

(<http://play.google.com/store/apps/details?id=com.analyticsvidhya.android>)

