

[f](https://www.facebook.com/AnalyticsVidhya) (<https://www.facebook.com/AnalyticsVidhya>) | [t](https://twitter.com/analyticsvidhya) (<https://twitter.com/analyticsvidhya>)

[g+](https://plus.google.com/+Analyticsvidhya/posts) (<https://plus.google.com/+Analyticsvidhya/posts>)

[in](https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165) (<https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165>)

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

Jobs (<https://www.analyticsvidhya.com/jobs/>) | Trainings (<https://trainings.analyticsvidhya.com>)

Learning Paths (<https://www.analyticsvidhya.com/learning-paths-data-science-business-analytics-business-intelligence-big-data/>)

Discuss (<https://discuss.analyticsvidhya.com>) | Corporate (<https://www.analyticsvidhya.com/corporate/>)



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



([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVhome\\_top](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVhome_top))

Home (<https://www.analyticsvidhya.com/>) > Machine Learning (<https://www.analyticsvidhya.com/blog/category/machine-learning/>) > 40 Interview Questions asked at Startups in Machine Learning / Data Science (<https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/>)

# 40 Interview Questions asked at Startups in Machine Learning / Data Science

MACHINE LEARNING (<https://www.analyticsvidhya.com/blog/category/machine-learning/>)

## TOP ANALYTICS VIDHYA USERS

Rank	Name
1	vopani ( <a href="https://datahack.analyticsvidhya.com/user/profile/Rohit">https://datahack.analyticsvidhya.com/user/profile/Rohit</a> )
2	SRK ( <a href="https://datahack.analyticsvidhya.com/user/profile/SHRINIVAS">https://datahack.analyticsvidhya.com/user/profile/SHRINIVAS</a> )
3	binga ( <a href="https://datahack.analyticsvidhya.com/user/profile/binga">https://datahack.analyticsvidhya.com/user/profile/binga</a> )

<https://www.facebook.com/sharer.php?u=https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science>  
<https://twitter.com/intent/tweet?url=https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science>

OScience+https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-

-machine-learning-data-science/) 8+ (https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-4

) P (https://www.analyticsvidhya.com) GET HIRED ns.png& 20in%20Machine%20Learning%20

iew-questions-asked-at-startups-in-machine-learning-data-science/&media=https://www.analyticsvidhya.com/blog/2016

ns.png& 20in%20Machine%20Learning%20

AVBYTES (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/)

Data Engineering

CONTACT US (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

Hackathon

WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

CAMPUS AMBASSADOR (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/)

/data-engineering-talent-hunt-hackathon

Carefully! These question can make you think

CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/)

**THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST**

By IEOR @ IITB & McKinsey  
Knowledge Centre India

REGISTER

Win Prizes Worth  
INR 1 Lakh+

(https://youroptimizationcontest.com/contests/the-young-optimization-crackerjack-contest)

/?utm\_source=AV&utm\_medium=Display&

However, getting into these roles is not easy. You obviously

need to get excited about the idea, team and the vision of the company. You might also find some real difficult technical questions on your way. The set of questions asked depend on what does the startup do. Do they provide consulting? Do they build ML products? You should always find this out prior to beginning your interview preparation.

To help you prepare for your next interview, I've prepared a list of 40 plausible & tricky questions which are likely to come across your way in interviews. If you can answer and understand these questions, rest assured, you will give a tough fight in your job interview.

Aayushmnit  
LEARN ENGAGE COMPETE  
4 (https://datahack.  
/user/profile/aay

Mark Landry  
5 (https://datahack.  
/user/profile/mark

More Rankings

(http://datahack.analyticsvidhya.com  
/users)

Answer Questions as Fast as You Can Think



LEADERS IN ANALYTICS

(http://www.greatlearning.edu  
/analytics  
/?utm\_source=avm&  
utm\_medium=avmbanner&  
utm\_campaign=pgpba+bda)



(https://upgrad.com/data-  
science?utm\_source=AV&  
utm\_medium=Display&  
utm\_campaign=DS\_AV\_Banner  
utm\_term=DS\_AV\_Banner&  
utm\_content=DS\_AV\_Banner)

Note: A key to answer these questions is to have concrete practical understanding on ML and related statistical concepts.

<https://www.analyticsvidhya.com>

LEARN > ENGAGE > COMPETE >



Asked at <https://www.analyticsvidhya.com/campus-ambassador/>

/data-engineering-talent-hunt-hackathon

ng-talent-hunt-hackathon

CORPORATE ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://www.analyticsvidhya.com/corporate/))

# THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

**REGISTER**

**Win Prizes Worth  
INR 1 Lakhs**

# Machine Learning

aving 1000 columns and 1  
ased on a classification  
ked you to reduce the  
odel computation time can  
memory constraints. What  
e practical assumptions )

**Answer:** Processing a high dimensional data on a limited memory machine is a strenuous task, your interviewer would be fully aware of that. Following are the methods you can use to tackle such situation:

1. Since we have lower RAM, we should close all other applications in our machine, including the web browser, so that most of the memory can be put to use.
  2. We can randomly sample the data set. This means, we can create a smaller data set, let's say, having 1000 variables and 300000 rows and do the computations.
  3. To reduce dimensionality, we can separate the numerical and

## POPULAR POSTS

- Essentials of Machine Learning Algorithms (with Python and R Codes)  
(<https://www.analyticsvidhya.com/blog/2017/09/common-machine-learning-algorithms/>)
  - A Complete Tutorial to Learn Data Science with Python from Scratch  
(<https://www.analyticsvidhya.com/blog/2016/01/complete-tutorial-learn-data-science-python-scratch-2/>)
  - 7 Types of Regression Techniques you should know!  
(<https://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/>)
  - 11 most read Deep Learning Articles from

categorical variables and remove the correlated variables. For numerical variables, we'll use correlation. For categorical variables, we'll use chi-square test.

4. Also, we can use PCA (<https://www.analyticsvidhya.com/blog/2016/09/practical-guide-principal-component-analysis/>)



nponents which can explain the variance in the data set.

ake Vowpal Wabbit (available <https://www.vowpalwabbit.org/>)

chastic Gradient Descent is

understanding to estimate the response variable. But,

this is an iterative approach failing to identify useful predictors might result in significant loss of information.

[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](https://www.analyticsvidhya.com/corporate/)

read about online learning (<https://www.analyticsvidhya.com/blog/2015/09/online-learning-simplified-2/>) & (<https://www.quora.com/What-are-the-differences-between-gradient-descent-and-stochastic-gradient-descent-methods>).

[THE YOUNG OPTIMIZATION CRACKERJACK CONTEST](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/closing-ceremony/)

By IEOR @ IITB & McKinsey Knowledge Centre India

[Win Prizes Worth INR 1 Lakh+](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/closing-ceremony/)

[DATAFEST 2017](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/closing-ceremony/)

[Register Now!](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/closing-ceremony/)

(<https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/closing-ceremony/>)

?utm\_source=AVblog\_sidebottom )

## Q2. Is rotation necessary in PCA? If yes, Why? What will happen if you don't rotate the components?

**Answer:** Yes, rotation (orthogonal) is necessary because it maximizes the difference between variance captured by the component. This makes the components easier to interpret. Not to forget, that's the motive of doing PCA where, we aim to select fewer components (than features) which can explain the maximum variance in the data set. By doing rotation, the relative location of the components doesn't change, it only changes the actual coordinates of the points.

Analytics Vidhya in 2017  
LEARN ✓ ENGAGE ✓ COMPETE ✓  
<https://www.analyticsvidhya.com/blog/2017/12/11-deep-learning-analytics-vidhya-2017/>

Understanding Support Vector Machine algorithm from examples (along with code)

(<https://www.analyticsvidhya.com/blog/2017/09/understanding-support-vector-machine-example-code/>)

A Complete Tutorial on Tree Based Modeling from Scratch (in R & Python)

(<https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/>)

10 Data Science, Machine Learning and AI Podcasts You Must Listen To

(<https://www.analyticsvidhya.com/blog/2018/01/10-data-science-machine-learning-ai-podcasts-must-listen/>)

6 Easy Steps to Learn Naive Bayes Algorithm (with codes in Python and R)

(<https://www.analyticsvidhya.com/blog/2017/09/naive-bayes-explained/>)

If we don't rotate the components, the effect of PCA will diminish. [HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/) to select more number of components to explain variance in the data set.

LEARN ▾ ENGAGE ▾ COMPETE ▾ [SEARCH](#)

<https://www.analyticsvidhya.com/> GET HIRED ▾



[icsvidhya.com/blog/2016](https://icsvidhya.com/blog/2016)

ment analysis python/)

ta set has missing values

dition from the median

ain unaffected? Why?

(<https://www.analyticsvidhya.com/campus-ambassador/>)

**Answer:** This question has enough hints for you to start /data-engineering-talent-hunt-hackathon

thinking. Since the data is spread across median, let's assume

CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/), in a normal distribution, standard deviation from mean (or 68% of the data unaffected). Median would remain unaffected by

the data engineering talent hunt hackathon.



([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblog\\_sidebottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblog_sidebottom))

**Answer:** If you have worked on enough data sets, you should deduce that cancer detection results in imbalanced data. In an imbalanced data set, accuracy should not be used as a measure of performance because 96% (as given) might only be predicting majority class correctly, but our class of interest is minority class (4%) which is the people who actually got diagnosed with cancer. Hence, in order to evaluate model performance, we should use Sensitivity (True Positive Rate), Specificity (True Negative Rate), F measure to determine class wise performance of the classifier. If the minority class performance is found to be poor, we can undertake the

## RECENT POSTS



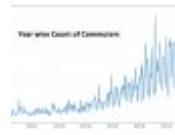
TextBlob

(<https://www.analyticsvidhya.com/blog/2018/02/natural-language-processing-for-beginners-using-textblob/>)

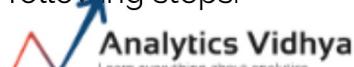
Natural Language Processing for Beginners: Using TextBlob

(<https://www.analyticsvidhya.com/blog/2018/02/natural-language-processing-for-beginners-using-textblob/>)

SHUBHAM JAIN , FE...



following steps:



1. We can use undersampling, oversampling or SMOTE to make the data balanced.

[GET HIRED](#)



(<https://www.analyticsvidhya.com/campus-ambassador/>)

(<https://www.analyticsvidhya.com/blog/2016/03/practical-data-engineering-talent-hunt-hackathon-guide-dealing-imbalanced-classification-problems/>)

[CORPORATE \(<https://www.analyticsvidhya.com/corporate/>\)](https://www.analyticsvidhya.com/corporate/)

## THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey Knowledge Centre India

[REGISTER](#)

Win Prizes Worth  
INR 1 Lakh+

(<https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest>)

**Q6. Explain prior probability, likelihood and marginal likelihood in context of naiveBayes algorithm?**

**Answer:** Prior probability is nothing but, the proportion of dependent (binary) variable in the data set. It is the closest guess you can make about a class, without any further information. For example: In a data set, the dependent variable is binary (1 and 0). The proportion of 1 (spam) is 70% and 0 (not spam) is 30%. Hence, we can estimate that there are 70% chances that any new email would be classified as spam.

Likelihood is the probability of classifying a given observation as 1 in presence of some other variable. For example: The

(<https://www.analyticsvidhya.com/blog/2018/02/time-series-forecasting-methods/>)

**7 methods to perform Time Series forecasting (with Python codes)**  
(<https://www.analyticsvidhya.com/blog/2018/02/time-series-forecasting-methods/>)

GURCHETAN SINGH ...



(<https://www.analyticsvidhya.com/blog/2018/02/introductory-guide-regularized-greedy-forests-rgf-python/>)

**An Introductory Guide to Regularized Greedy Forests (RGF) with a case study in Python**  
(<https://www.analyticsvidhya.com/blog/2018/02/introductory-guide-regularized-greedy-forests-rgf-python/>)

ANKIT CHOUDHARY ...



probability that the word 'FREE' is used in previous spam messages like [HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/). Marginal likelihood is, the probability that the word 'FREE' is used in any message.

<https://www.analyticsvidhya.com/> GET HIRED ▾



(<https://www.analyticsvidhya.com/blog/2018/02/demystifying-security-data-science/>)

LEARN ▾ ENGAGE ▾ COMPETE ▾ (<https://www.analyticsvidhya.com/blog/2018/02/demystifying-security-data-science/>)

Demystifying Information Security Using Data Science (<https://www.analyticsvidhya.com/blog/2018/02/demystifying-data-science/>)

GUEST BLOG , FEBR...

**Answer:** Time series data is known to posses linearity. On the

[CAMPUS AMBASSADOR \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/\)](https://www.analyticsvidhya.com/campus-ambassador/) Answer: Time series data is known to posses linearity. On the other hand, non-linear data is known to work best to predict future values. The reason why decision tree models fail to predict future values is because it couldn't map the non-linear relationship which a regression model did. A linear regression model can only predict future values if the data set satisfies its linearity assumption. (<https://www.analyticsvidhya.com/blog/2016/09/assumptions-plots/>)

([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblogidebottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblogidebottom))

**Q8. You are assigned a new project which involves helping a food delivery company save more money. The problem is, company's delivery team aren't able to deliver food on time. As a result, their customers get unhappy. And, to keep them happy, they end up delivering food for free. Which machine learning algorithm can save them?**

**Answer:** You might have started hopping through the list of ML algorithms in your mind. But, wait! Such questions are asked to test your machine learning fundamentals.

This is not a machine learning problem. This is a route optimization problem. A machine learning problem consist of



([http://www.edvancer.in/certified-data-scientist-with-python-course?utm\\_source=AV&utm\\_medium=AVads&](http://www.edvancer.in/certified-data-scientist-with-python-course?utm_source=AV&utm_medium=AVads&))

three things:



1. There exist a pattern.  
(<https://www.analyticsvidhya.com>)

2. You cannot solve it mathematically (even by writing

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)

utm\_campaign=AVadsnonfc&utm\_content=pythondavad) [LEARN](#) [ENGAGE](#) [COMPETE](#) [SEARCH](#)

[GET HIRED](#)



[LOW BIAS AND HIGH VARIANCE, WHICH ALGORITHM SHOULD YOU USE TO TACKLE IT? WHY?](https://campus.ambassadorscity.com/analyticsvidhya.com/campus-ambassador-to-tackle-it-why?utm_source=AVblog&utm_medium=Sidebottom)

([https://campus.ambassadorscity.com/analyticsvidhya.com/campus-ambassador-to-tackle-it-why?utm\\_source=AVblog&utm\\_medium=Sidebottom](https://campus.ambassadorscity.com/analyticsvidhya.com/campus-ambassador-to-tackle-it-why?utm_source=AVblog&utm_medium=Sidebottom))



By IEOR @ IITB & McKinsey Knowledge Centre India

[REGISTER](#)

Win Prizes Worth  
INR 1 Lakh+

the model's predicted values words, the model becomes data distribution. While it but not to forget, a flexible abilities. It means, when this data, it gives disappointing

to resubmit to tackle high variance problem. Bagging algorithms divides a data set into subsets made with repeated randomized Sampling. Then, these samples are used to generate a set of models using a single learning algorithm. Later, the model predictions are combined using voting (classification) or averaging (regression).

Also, to combat high variance, we can:

1. Use regularization technique, where higher model coefficients get penalized, hence lowering model complexity.
2. Use top n features from variable importance chart. May be, with all the variable in the data set, the algorithm is having difficulty in finding the meaningful signal.

## GET CONNECTED



14,247

FOLLOWERS

(<http://www.twitter.com/analyticsvidhya>)



42,848

FOLLOWERS

(<http://www.facebook.com/AnalyticsVidhya>)



2,550

FOLLOWERS

(<https://plus.google.com/+AnalyticsVidhya>)



Email

SUBSCRIBE

(<http://feedburner.google.com/fb/a/mailverify?uri=analyticsvidhya>)

**Q10.** You are given a data set. The data set contains many variables, some of which are highly correlated and you know about it. Your manager has asked you to run PCA. Would you remove correlated variables? Why?



I would prompt to say No, but that related variables have a presence of correlated a particular component

<https://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/>

a data set, of which 2 are

<https://CAMPUS-AMBASSADOR.INTERFACE.WEB.COM/ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/>  
/data-engineering-talent-hunt-hackathon

CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/). Also, adding correlated importance on those variable,



, you are now anxious to

a result, you build 5 GBM algorithm would do the magic.

could perform better than

[https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AnalyticsVidhya](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AnalyticsVidhya)  
benchmark score. Finally, you decided to combine those models. Though ensembled models are known to return high accuracy but are unfortunate. Where did you miss?

**Answer:** As we know, ensemble learners are based on the idea of combining weak learners to create strong learners. But, these learners provide superior result when the combined models are uncorrelated. Since, we have used 5 GBM models and got no accuracy improvement, suggests that the models are correlated. The problem with correlated models is, all the models provide same information.

For example: If model 1 has classified User1122 as 1, there are high chances model 2 and model 3 would have done the same, even if its actual value is 0. Therefore, ensemble learners are

built on the premise of combining weak uncorrelated models to obtain better predictions.

 Learn everything about analytics

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)

LEARN ▾ ENGAGE ▾ COMPETE ▾



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

GET HIRED ▾



ns clustering?

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](https://www.analyticsvidhya.com/blog/category/avbytes/)

**Data Engineering** their names. You should  
[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/) see between both these

**Hackathon**

WE ARE HIRING! (<https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/>)

(<https://www.analyticsvidhya.com/campus-ambassador/>)

kmeans algorithm partitions a data set into clusters such that a  
 /data-engineering-talent-hunt-hackathon

cluster formed is homogeneous and the points in each cluster

**CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/)** tries to maintain enough  
 members. Due to unsupervised

labeled observation based  
 finding neighbors. It is also  
 involves minimal training of  
 training data to make



REGISTER

Win Prizes Worth  
INR 1 Lakh+

(<https://datahack.analyticsvidhya.com/contest>

/the-young-optimization-crackerjack-contest

?utm\_source=AVblog\_sidebottom )

**Q13. How is True Positive Rate and Recall related? Write the equation.**

**Answer:** True Positive Rate = Recall. Yes, they are equal having the formula  $(TP / (TP + FN))$ .

Know more: Evaluation Metrics

(<https://www.analyticsvidhya.com/blog/2016/02/7-important-model-evaluation-error-metrics/>)

**Q14. You have built a multiple regression model. Your model**

$R^2$  isn't as good as you wanted. For improvement, your model  $R^2$  becomes 0.8 from removing the intercept term.

Learn everything about analytics  
0.3 Is it possible? How?  
(<https://www.analyticsvidhya.com>)  
GET HIRED

LEARN < ENGAGE < COMPETE < 



need to understand the  
in a regression  
**OM/BLOG/CATEGORY/AVBYTES/**)  
odel prediction without any  
tion. The formula of  $R^2 = 1 - \frac{\text{predicted value}}{\text{actual value}}$

WE ARE HIRING! (<https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/>)

of intercept term ( vmean )



our manager has informed  
ring from multicollinearity.  
true? Without losing any  
tter model?

we can create a correlation

matrix to identify & remove variables having correlation above 75% (deciding a threshold is subjective). In addition, we can use calculate VIF (variance inflation factor) to check the presence of multicollinearity.

multicollinearity. VIF value  $\leq 4$  suggests no multicollinearity whereas a value of  $\geq 10$  implies serious multicollinearity. Also, we can use tolerance as an indicator of multicollinearity.

But, removing correlated variables might lead to loss of information. In order to retain those variables, we can use penalized regression models like ridge or lasso regression. Also, we can add some random noise in correlated variable so that the variables become different from each other. But, adding noise might affect the prediction accuracy, hence this approach should be carefully used.

Know more: Regression (<https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions/#q17>)

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions/#q17)

LEARN ▾ ENGAGE ▾ COMPETE ▾



[plots-solutions/](https://www.analyticsvidhya.com/plots-solutions/)

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

GET HIRED ▾



AVBYTES ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/](https://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/))

CONTACT US ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://WWW.ANALYTICSVIDHYA.COM/CONTACT/))

Hackathon

WE ARE HIRING! ([HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/](https://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/))

In presence of many

variables with small to medium sized effect, use ridge

([https://CAMPUS-AMBASSADOR \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/\)](https://CAMPUS-AMBASSADOR (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/)))

/data-engineering-talent-hunt-hackathon

CORPORATE ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://WWW.ANALYTICSVIDHYA.COM/CORPORATE/))

regression (L1) does both shrinkage, whereas Ridge shrinkage and end up the model. In presence of correlation might be the preferred Lasso best in situations where higher variance. Therefore, it

THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

Win Prizes Worth  
INR 1 Lakh+

Lasso      Regression

REGISTER

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

(<https://www.analyticsvidhya.com/blog/2016/01/complete-the-young-optimization-crackerjack-contest-tutorial-ridge-lasso-regression-python/>)

**Q17. Rise in global average temperature led to decrease in number of pirates around the world. Does that mean that decrease in number of pirates caused the climate change?**

**Answer:** After reading this question, you should have understood that this is a classic case of "causation and correlation". No, we can't conclude that decrease in number of pirates caused the climate change because there might be other factors (lurking or confounding variables) influencing this phenomenon.

Therefore, there might be a correlation between global average temperature and number of pirates, but based on this information we can't say that pirates died because of rise in global average temperature.

LEARN ▾ ENGAGE ▾ COMPETE ▾ 



[\(https://www.analyticsvidhya.com/campus-ambassador/\)](https://www.analyticsvidhya.com/campus-ambassador/)

/data-engineering-talent-hunt-hackathon

Following are the methods of variable selection you

[CORPORATE \(https://www.analyticsvidhya.com/corporate/\)](https://www.analyticsvidhya.com/corporate/)

prior to selecting important

variables based on p values  
Forward Selection, Stepwise

and plot variable importance

By IEOR @ IITB & McKinsey

Knowledge Centre India

[REGISTER](#)

Win Prizes Worth  
INR 1 Lakh+

Measure information gain for the available set of features and

select top n features accordingly.

(https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\_source=AVblog\_sidebottom )

**Q19. What is the difference between covariance and correlation?**

**Answer:** Correlation is the standardized form of covariance.

Covariances are difficult to compare. For example: if we calculate the covariances of salary (\$) and age (years), we'll get different covariances which can't be compared because of having unequal scales. To combat such situation, we calculate correlation to get a value between -1 and 1, irrespective of their respective scale.



HOME ([HTTPS://WWW.ANALYTICSVIDHYA.COM/](https://WWW.ANALYTICSVIDHYA.COM/)) LEARN ENGAGE COMPETE

**Q20. Is it possible capture the correlation between continuous and categorical variable? If yes, how?**



([HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/](https://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/))

/data-engineering-talent-hunt-hackathon

**Answer:** The fundamental difference is, random forest uses

boosting ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://WWW.ANALYTICSVIDHYA.COM/CORPORATE/))

**THE YOUNG OPTIMIZATION CRACKERJACK CONTEST**

By IEOR @ IITB & McKinsey Knowledge Centre India

REGISTER

Win Prizes Worth INR 1 Lakh+

Answer: Random forest uses boosting. It divides the data into n samples using bagging. The single learning algorithm a tree is applied on each sample. The resultant predictions are combined using averaging. Bagging is done is sequential manner. In each round of predictions, the weights of misclassified predictions are increased. Such that they can be corrected in the succeeding round. This sequential process of giving higher weights to misclassified predictions continues until a stopping criterion is reached.

([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/THE-YOUNG-OPTIMIZATION-CRACKERJACK-CONTEST/?UTM\\_SOURCE=AVBLOG\\_SIDEBOTTOM](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblog_sidebottom))

Random forest improves model accuracy by reducing variance (mainly). The trees grown are uncorrelated to maximize the decrease in variance. On the other hand, GBM improves accuracy by reducing both bias and variance in a model.

Know more: Tree based modeling

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/04/COMPLETE-TUTORIAL-TREE-BASED-MODELING-SCRATCH-IN-PYTHON/](https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/))

**Q22. Running a binary classification tree algorithm is the easy**

part. Do you know how does a tree splitting takes place i.e. how does this field decide which variable to split at the root node and succeeding nodes? )

HOME (HTTPS://WWW.ANALYTICSVIDHYA.COM/) LEARN ENGAGE COMPETE



This decision based on Gini

words, the tree algorithm can divide the data set into

CONTACT US (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

Hackathon

WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

ms from a population at

described by a probability distribution.

calculate Gini as following:

(https://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/)

1. Calculate Gini for sub-nodes, using formula sum of square of probability for success and failure ( $p^2+q^2$ ).



Integrated Gini score of each node

is given by (for binary class):

$$= p \log_2 p + q \log_2 q$$

$p$  and  $q$  are the probabilities of success and failure respectively

If all the samples in a node belong to one class, then node is homogeneous. It is

present in a node at 50% probability.

(https://datahackathon.analyticsvidhya.com/contest

/the-young-optimization-crackerjack-contest

?utm\_source=AVblog\_sidebottom )

**Q23. You've built a random forest model with 10000 trees. You got delighted after getting training error as 0.00. But, the validation error is 34.23. What is going on? Haven't you trained your model perfectly?**

**Answer:** The model has overfitted. Training error 0.00 means the classifier has mimiced the training data patterns to an extent, that they are not available in the unseen data. Hence, when this classifier was run on unseen sample, it couldn't find those patterns and returned prediction with higher error. In random forest, it happens when we use larger number of trees

than necessary. Hence, to avoid these situations, we should tune number of variables used across validation.

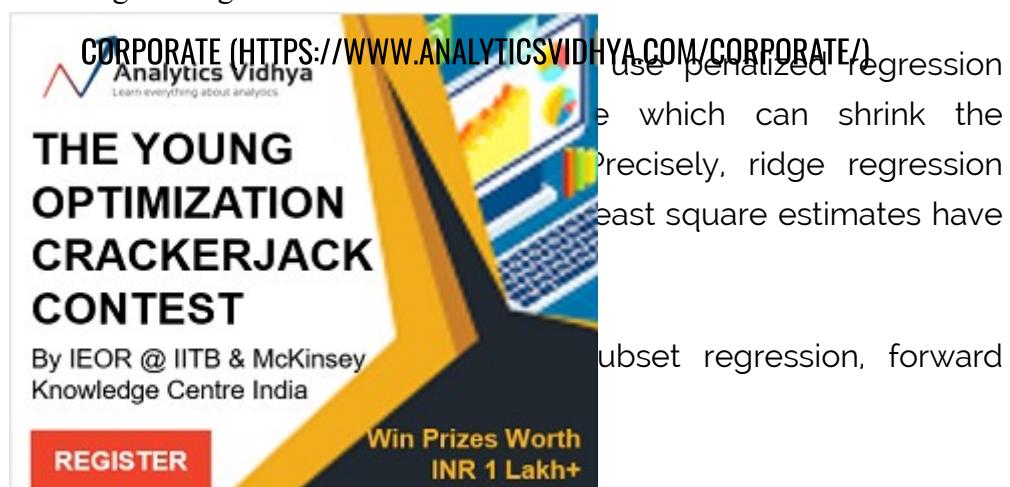
HOME (HTTPS://WWW.ANALYTICSVIDHYA.COM/) LEARN ▾ ENGAGE ▾ COMPETE ▾ 

(https://www.analyticsvidhya.com )  
GET HIRED ▾



WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)  
(https://CAMPUS-AMBASSADOR (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/)

O/S cannot be used at all  
/data-engineering-talent-hunt-hackathon

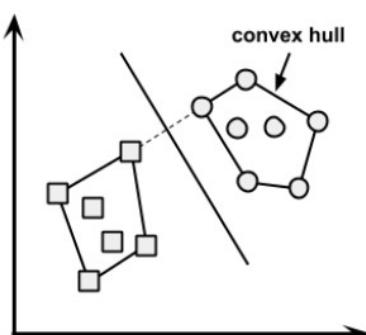


By IEOR @ IITB & McKinsey Knowledge Centre India  
Win Prizes Worth INR 1 Lakh+  
[REGISTER](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest?utm_source=AVblog_sidebottom)

subset regression, forward

(https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest?utm\_source=AVblog\_sidebottom )  
Q25. What is convex hull ?  
(Hint: Think SVM)

**Answer:** In case of linearly separable data, convex hull represents the outer boundaries of the two groups of data points. Once convex hull is created, we get maximum margin hyperplane (MMH) as a perpendicular bisector between two convex hulls. MMH is the line which attempts to create greatest separation between two groups.



**Q26. We know that one hot encoding increasing the dimensionality of our data set. But, label encoding doesn't.**

How? (<https://www.analyticsvidhya.com>)  
GET HIRED ▾



tion. It's a simple question

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

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

Hackathon

WE ARE HIRING! (<https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/>)

coding 'color' variable will

(<https://campus.ambassador.analyticsvidhya.com/campus-ambassador/>)

/dataColor. Create a new categorical variable containing 0 and 1 value.



ue would you use on time  
V?

(<https://datahack.analyticsvidhya.com/contest>

/the-young-optimization-crackerjack-contest

In time series problem, k fold can be troublesome because  
/?utm\_source=AVblog\_sidebottom there might be some pattern in year 4 or 5 which is not in year

3. Resampling the data set will separate these trends, and we might end up validation on past years, which is incorrect. Instead, we can use forward chaining strategy with 5 fold as shown below:

- fold 1 : training [1], test [2]
- fold 2 : training [1 2], test [3]
- fold 3 : training [1 2 3], test [4]
- fold 4 : training [1 2 3 4], test [5]
- fold 5 : training [1 2 3 4 5], test [6]

where 1,2,3,4,5,6 represents "year".



HOME ([HTTPS://WWW.ANALYTICSVIDHYA.COM/](https://WWW.ANALYTICSVIDHYA.COM/)) LEARN ENGAGE COMPETE

**Q28. You are given a data set consisting of variables having more than 30% missing values.** Let's say, out of 50 variables,



higher than 30%. How will

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](AVBYTES (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/))

**Data Engineering**

**Hackathon**

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/))

distribution with the target

[CAMPUS AMBASSADOR \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/\)](CAMPUS AMBASSADOR (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/))  
 variable, and if found any pattern we'll keep those missing  
 /data-engineering-talent-hunt-hackathon values and assign them a new category while removing

[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/))

**THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST**

By IEOR @ IITB & McKinsey  
Knowledge Centre India

**REGISTER**

Win Prizes Worth  
INR 1 Lakh+

nsiders "User Behavior" for

<RECOMMENDING ITEMS! They exploit behavior of other users and the-young-optimization-crackerjack-contest-items-in-terms-of-transaction-history,-ratings,-selection-and-purchase-information. Other users behaviour and preferences over the items are used to recommend items to the new users. In this case, features of the items are not known.>

Know more: Recommender System

<https://www.analyticsvidhya.com/blog/2015/10/recommendation-engines/>

**Q30. What do you understand by Type I vs Type II error ?**

**Answer:** Type I error is committed when the null hypothesis is

true and we reject it, also known as a 'False Positive'. Type II error is committed when the null hypothesis is false and we accept it, also known as 'False Negative'.

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](#) LEARN ▾ ENGAGE ▾ COMPETE ▾ [GET HIRED](#) ▾ [SEARCH](#)



we can say Type I error is positive (1) when it is occurs when we classify a

[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](#)

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](#)

classification problem. For validation purposes, you've randomly sampled the training

([https://CAMPUS-AMBASSADOR.IHCSVIDHYA.COM/ANALYTICSVIDHYA.COM/CAMPOS-AMBASSADOR/](#))  
/data-engineering-talent-hunt-hackathon



n unseen data since your ever, you get shocked after ent wrong?

problem, we should always random sampling. A random consideration the proportion of stratified sampling helps to variable in the resultant

([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblog\\_sidebottom](#) )

**Q32. You have been asked to evaluate a regression model based on R<sup>2</sup>, adjusted R<sup>2</sup> and tolerance. What will be your criteria?**

**Answer:** Tolerance (1 / VIF) is used as an indicator of multicollinearity. It is an indicator of percent of variance in a predictor which cannot be accounted by other predictors. Large values of tolerance is desirable.

We will consider adjusted R<sup>2</sup> as opposed to R<sup>2</sup> to evaluate model fit because R<sup>2</sup> increases irrespective of improvement in prediction accuracy as we add more variables. But, adjusted R<sup>2</sup> would only increase if an additional variable improves the

accuracy of model, otherwise stays same. It is difficult to commit threshold value for adjusted R<sup>2</sup> because it varies between data sets. For example: a gene mutation data set might result in lower adjusted R<sup>2</sup> and still provide fairly good

LEARN ▾ ENGAGE ▾ COMPETE ▾ 



(<https://www.analyticsvidhya.com/campus-ambassador/>)

**Answer:** We don't use manhattan distance because

or vertically only. It has

hand, euclidean metric can

distance. Since, the data

ension, euclidean distance is a

the movement made by a  
anhattan distance because  
tal movements.

(<https://datahack.analyticsvidhya.com/contest>

/the-young-optimization-crackerjack-contest

**Q34. Explain machine learning to me like a 5 year old.**

?utm\_source=AVblog\_sidebottom )

**Answer:** It's simple. It's just like how babies learn to walk. Every time they fall down, they learn (unconsciously) & realize that their legs should be straight and not in a bend position. The next time they fall down, they feel pain. They cry. But, they learn 'not to stand like that again'. In order to avoid that pain, they try harder. To succeed, they even seek support from the door or wall or anything near them, which helps them stand firm.

This is how a machine works & develops intuition from its environment.

Note: The interview is only trying to test if have the ability of

explain complex concepts in simple terms.



[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](#) LEARN ▾ ENGAGE ▾ COMPETE ▾

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

[GET HIRED ▾](#)



tion model is generally

value. How would you

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](#)

? [CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](#)

[DATA ENGINEERING Hacks](#)

[Hackathon](#)

WE ARE HIRING! (<https://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/>)

to predict probabilities we

with confusion matrix to

(<https://www.analyticsvidhya.com/campus-ambassador>)

/data-engineering-talent-hunt-hackathon

regarding the analogous methods of selecting the best

regression models. AIC is the measure of fit which penalizes

the number of coefficients. Therefore, we

try to minimize the AIC value.

The AIC is the measure of fit of a response predicted by a model

to its data. Lower the value, better the fit.

AIC = -2 ln(L) + 2k where L denotes the response predicted by the model on k independent variables. Lower the

AIC value, better the fit of the model.

THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey Knowledge Centre India

Win Prizes Worth INR 1 Lakh+

[https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblog\\_sidebottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblog_sidebottom)

blog/2015/11/beginners-

contests-for-data-science-students/

([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblog\\_sidebottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblog_sidebottom) )

**Q36. Considering the long list of machine learning algorithm, given a data set, how do you decide which one to use?**

**Answer:** You should say, the choice of machine learning algorithm solely depends of the type of data. If you are given a data set which is exhibits linearity, then linear regression would be the best algorithm to use. If you given to work on images, audios, then neural network would help you to build a robust model.

If the data comprises of non linear interactions, then a boosting or bagging algorithm should be the choice. If the business

requirement is to build a model which can be deployed, then we'll use **Analytics Vidhya** ([HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)) a decision tree model (easy to interpret and explain) instead of black box algorithms like SVM, GBM etc.

LEARN ▾ ENGAGE ▾ COMPETE ▾ 

[GET HIRED ▾](#)



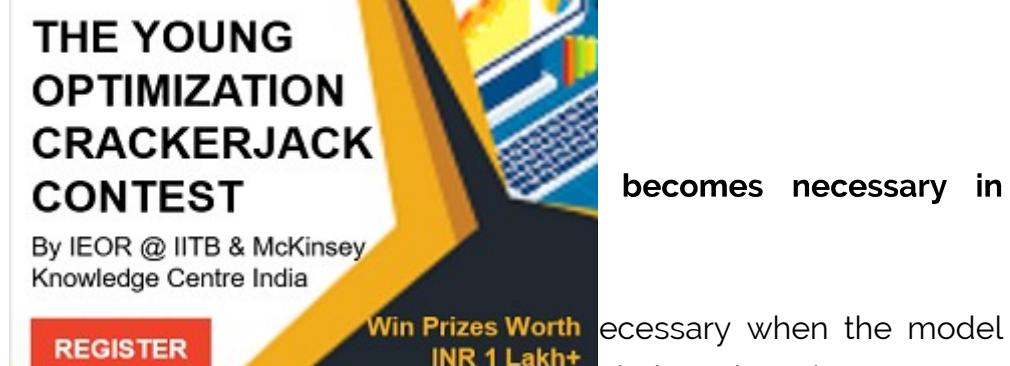
(<https://www.analyticsvidhya.com/campus-ambassador/>)

/data-engineering-talent-hunt-hackathon

**CORPORATE** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://WWW.ANALYTICSVIDHYA.COM/CORPORATE/))

categorical variable can be

only when the variable is



ecessary when the model chnique introduces a cost

(<https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest>)

term for bringing in more features with the objective function.

/the-young-optimization-crackerjack-contest

Hence, it tries to push the coefficients for many variables to

?utm\_source=AVblog\_sidebottom

zero and hence reduce cost term. This helps to reduce model

complexity so that the model can become better at predicting

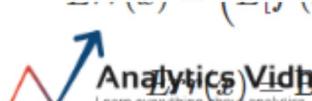
(generalizing).

becomes necessary in

### Q39. What do you understand by Bias Variance trade off?

**Answer:** The error emerging from any model can be broken down into three components mathematically. Following are these component :

$$Err(x) = \left( E[\hat{f}(x)] - f(x) \right)^2 + E\left[ \hat{f}(x) - E[\hat{f}(x)] \right]^2 + \sigma^2$$



[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)

LEARN ▾ ENGAGE ▾ COMPETE ▾



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

(<https://www.analyticsvidhya.com/wp-content/uploads>

**ABInBev**

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](https://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/)

such on an average are the

[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

actual value. A high bias

**Hackathon**

learning model which keeps

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](https://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

on the other side made on same observation

bias. A low bias

model will over-fit on

(<https://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/>)

/data-engineering-talent-hunt-hackathon

[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](https://WWW.ANALYTICSVIDHYA.COM/CORPORATE/)

**CORPORATE**  
Analytics Vidhya  
Learn everything about analytics

**THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST**

By IEOR @ IITB & McKinsey  
Knowledge Centre India

**REGISTER**

Win Prizes Worth  
INR 1 Lakh+

**Maximum likelihood is to**  
choose the best parameter.

Good are the methods used  
to approximate the  
value. In simple words,

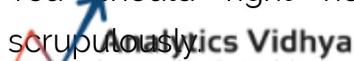
([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVBlog&utm\\_medium=9](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVBlog&utm_medium=9))

Ordinary least square(OLS) is a method used in linear regression which approximates the parameters resulting in minimum distance between actual and predicted values. Maximum Likelihood helps in choosing the values of parameters which maximizes the likelihood that the parameters are most likely to produce observed data.

## End Notes

You might have been able to answer all the questions, but the real value is in understanding them and generalizing your knowledge on similar questions. If you have struggled at these questions, no worries, now is the time to learn and not perform.

You should right now focus on learning these topics



[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)

LEARN ▾ ENGAGE ▾ COMPETE ▾



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

These questions are meant to give you a wide exposure on the

[GET HIRED](#)

in machine learning. I'm

curious enough to do

if you are planning for it.



[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](https://www.analyticsvidhya.com/contact/)

Have you appeared in any

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/)

I'd love to know your

(<https://www.analyticsvidhya.com/campus-ambassador/>)

/data-engineering-talent-hunt-hackathon



? Check out currently  
[www.analyticsvidhya.com](http://www.analyticsvidhya.com)  
1 data science.

[6/09/40-interview-questions-  
ence/?share=linkedin&nb=1\)](https://www.analyticsvidhya.com/corporate/?share=linkedin&nb=1)

1K+

[6/09/40-interview-questions-  
ence/?share=facebook&nb=1\)](https://www.analyticsvidhya.com/corporate/?share=facebook&nb=1)

666

[6/09/40-interview-questions-  
ence/?share=google-plus-1&nb=1\)](https://www.analyticsvidhya.com/corporate/?share=google-plus-1&nb=1)

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

/the-young-optimization-crackerjack-contest

([https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-  
asked-at-startups-in-machine-learning-data-science/?share=twitter&nb=1](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-<br/>asked-at-startups-in-machine-learning-data-science/?share=twitter&nb=1))

([https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-  
asked-at-startups-in-machine-learning-data-science/?share=pocket&nb=1](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-<br/>asked-at-startups-in-machine-learning-data-science/?share=pocket&nb=1))

([https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-  
asked-at-startups-in-machine-learning-data-science/?share=reddit&nb=1](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-<br/>asked-at-startups-in-machine-learning-data-science/?share=reddit&nb=1))

## RELATED

[Analytics Vidhya](#)



(https://www.analyticsvidhya.com/pitching-your-company-to-y-combinator)

HOME (HTTPS://WWW.ANALYTICSVIDHYA.COM/)

LEARN ▾ ENGAGE ▾ COMPETE ▾



GET HIRED



Gregory Piatetsky-Shapiro

Interview  
with  
Data Scientist

Ph.D. President, KDnuggets

Analytics Vidhya

ABInBev

AVBYTES (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/)

(https://www.analyticsvidhya.com/blog/2015

Data Engineering 16

CONTACT US (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

(https://www.analyticsvidhya.com/contact/)

Hackathon

WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

(https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/)

CAMPUS AMBASSADOR (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/)

(https://www.analyticsvidhya.com/campus-ambassador/)

Data Science Startups

Analytics Vidhya for

/data-engineering-talent-hunt-hackathon

from Y Combinator

the year 2016

CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/)

(https://www.analyticsvidhya.com/corporate/)

2016

Corporate

October 13, 2015

In "Business

Analytics"

THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

REGISTER

Win Prizes Worth  
INR 1 Lakh+

(https://datahack.analyticsvidhya.com/contest

/the-young-optimization-crackerjack-contest

?utm\_source=AVblog\_sidebottom\_ )

TAGS: BAGGING (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/BAGGING/), BIAS VARIANCE TRADEOFF (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/BIAS-VARIANCE-TRADEOFF/), BOOSTING (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/BOOSTING/), CROSS-VALIDATION (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/CROSS-VALIDATION/), DATA SCIENCE IN STARTUPS (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DATA-SCIENCE-IN-STARTUPS/), DATA SCIENTIST IN STARTUPS (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DATA-SCIENTIST-IN-STARTUPS/), DATA SCIENTIST INTERVIEW (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DATA-SCIENTIST-INTERVIEW/), INTERVIEW QUESTIONS (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/INTERVIEW-QUESTIONS/), MACHINE LEARNING (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MACHINE-LEARNING/), MACHINE LEARNING ENGINEER INTERVIEW QUESTION (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MACHINE-LEARNING-ENGINEER-INTERVIEW-QUESTION/), MAXIMUM LIKELIHOOD (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/MAXIMUM-LIKELIHOOD/), ORDINARY LEAST SQUARE (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/ORDINARY-LEAST-SQUARE/), REGULARIZATION (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/REGULARIZATION/)



[Previous Article](#) [HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](#) LEARN ENGAGE COMPETE [SEARCH](#)  
**AWS / Cloud Engineer –**  
 (https://www.analyticsvidhya.com )

Comprehensive

Introduction to Apache

Spark, RDDs &

Dataframes (using  
PySpark)

[Data Engineering](#) (https://www.analyticsvidhya.com)

[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/2016](#)

/09/comprehensive-

introduction-to-apache-

spark-rdds-dataframes-

using-pyspark/)

WE ARE HIRING! (https://www.analyticsvidhya.com/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

(https://www.analyticsvidhya.com/campus-ambassador/)

/data-engineering-talent-hunt-hackathon



[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](#)

## THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

[REGISTER](#)

Win Prizes Worth  
INR 1 Lakh+

Join our Content Team  
[analyticsvidhya.com  
contentteam/](#)

Content team

(https://datahack.analyticsvidhya.com/contest

/the-young-optimization-crackerjack-contest

?utm\_source=AVblog\_sidebottom )

This article is quiet old now and you might not get a prompt response from the author. We would request you to post this comment on Analytics Vidhya **Discussion portal** (<https://discuss.analyticsvidhya.com/>) to get your queries resolved.

## 33 COMMENTS

[REPLY \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-SCIENCE/COMMENTS/REPLYTOCOMMENT#RESPONSE\)](#)  
[www.analyticsvidhya.com](#)  
[/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-](#)

LEARNING-DATA-SCIENCE/#COMMENT-116110)



HOME (HTTPS://WWW.ANALYTICSVIDHYA.COM/)

LEARN ▾ ENGAGE ▾ COMPETE ▾



(https://www.analyticsvidhya.com )

GET HIRED ▾



content/Team says

INTERVIEW-QUESTIONS-

ASKED-AT-STARTUPS-IN-MACHINE-

-COMMENT-116122)

CONTACT US (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)

Hackathon

These questions help you to

passing interview rounds. All

WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

(https://www.analyticsvidhya.com/campus-ambassador/)

/data-engineering-talent-hunt-hackathon

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-

CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/COMPANY-CORPORATE/)

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-

-COMMENT-116111)

full to face on the true reality  
me 😊

## THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

REGISTER

Win Prizes Worth  
INR 1 Lakh+

content/Team says

INTERVIEW-QUESTIONS-

ASCIENCE/?REPLYTOCOM=116119#RESPOND)

VIDHYA.COM/BLOG/2016/09/40-

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-

(https://datahack.analyticsvidhya.com/young-data-science/#COMMENT-116119)

/the-young-optimization-crackerjack-contest

Hi Gianni

/?utm\_source=AVblog\_sidebottom )

Good to know, you found them helpful! All the

best.

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/COMMENT-116121#REPLYTOCOM=116121#RESPOND)

AnalyticsVidhya Content/Team says

(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-

LEARNING-DATA-SCIENCE/#COMMENT-116121)

Hi Gianni, I am happy to know that these  
question would help you in your journey. All the  
best.

REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116113](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#comment-116113))

LEARN ▾ ENGAGE ▾ COMPETE ▾ 

**GET HIRED** Good collection compiled by you Mr Manish ! Kudos !



CONTACT US ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))

**Hackathon**

WE ARE HIRING! ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CAREER/](https://www.analyticsvidhya.com/career/))

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116123#RESPOND

([HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/](https://www.analyticsvidhya.com/campus-ambassador/))

/data-engineering-talent-hunt-hackathon

Hi Prof Ravi You are right. These questions can

CORPORATE ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://www.analyticsvidhya.com/corporate/))

startups, facing off ML

question have higher

have laid emphasis on

as well.

## THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey  
Knowledge Centre India

REGISTER

Win Prizes Worth  
INR 1 Lakh+

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116117#RESPOND

([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST](https://datahack.analyticsvidhya.com/contest))

/the-young-optimization-crackerjack-contest

Thank you Manish Helpful for Beginners like me.

?utm\_source=AVblog\_sidebottom )

AnalyticsVidhya Content Team says INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116118#RESPOND

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116118](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#comment-116118))

Welcome 😊

chibote says INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116121#RESPOND

([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116124](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#comment-116124))



It seems Stastics is at the centre of Machine Learning.

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](#)

LEARN ▾ ENGAGE ▾ COMPETE ▾



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

[GET HIRED ▾](#)

[REPLY \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-#COMMENT-116126#RESPOND\)](#)

**ABInBev**

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](#)

[DATA SCIENCE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-#COMMENT-116126#RESPOND\)](#)

[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](#)

[HACKATHON](#)

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](#)

[INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-#COMMENT-116127#RESPOND\)](#)

(<https://campus.ambassador.analyticsvidhya.com>)

[/data-engineering-talent-hunt-hackathon](#)

[INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-#COMMENT-116127#RESPOND\)](#)

[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](#)

an inevitable part of machine learning. It is important to understand statistical concepts to master machine learning.

## THE YOUNG OPTIMIZATION CRACKERJACK CONTEST

By IEOR @ IITB & McKinsey Knowledge Centre India

[REGISTER](#)

[Win Prizes Worth INR 1 Lakh+ \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-#COMMENT-116128\)](#)

(<https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest>)

?utm\_source=AVblog\_sidebottom

I was wondering, do you recommend for somebody to special in a specific field of ML? I mean, it is recommended to choose between supervised learning and unsupervised learning algorithms, and simply say my specialty is this during an interview. Shouldn't organizations recruiting specify their specialty requirements too?

....and thank you for the post.



(https://www.analyticsvidhya.com/)

**AnalyticsVidhya Content**  
**HOME (HTTPS://WWW.ANALYTICSVIDHYA.COM/)** LEARN ▾ ENGAGE ▾ COMPETE ▾

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116129#RESPOND) IN-MACHINE-LEARNING-DATA-SCIENCE  
 (REPLYTOCOM=116129#RESPOND) (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116129#RESPOND)  
**GET HIRED** ▾ /BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE



WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)

(https://www.analyticsvidhya.com/campus-ambassador/)

/data-engineering-talent-hunt-hackathon

Hi Chibole,

It's always a good thing to establish yourself as an expert in a specific field. This helps the recruiter to understand that you are a detailed oriented person in machine learning, thinking of building your expertise in supervised learning.

would be good, but companies want more than that. Considering, the variety of data these days, they want someone who can deal with unlabeled data also. In short, they look for someone who isn't just an expert in operating Sniper Gun, but can use other weapons also if needed.



By IEOR @ IITB & McKinsey Knowledge Centre India

**REGISTER**

Win Prizes Worth  
INR 1 Lakh+

(https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\_source=AVblog\_sidebottom )

REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116125#RESPOND) IN-MACHINE-LEARNING-DATA-SCIENCE  
 (REPLYTOCOM=116125#RESPOND) (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116125)

\* stastics = Statistics



Karthikayan Sankaran (<https://www.linkedin.com/in/karthikayansankaran/>) says STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE  
 (REPLYTOCOM=116135#RESPOND) (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116135)



Hi Manish – Interesting & Informative set of questions &

answers. Thanks for compiling the same.

[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/) LEARN ▾ ENGAGE ▾ COMPETE ▾

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

GET HIRED ▾



**content/Team says**  
[INTERVIEW-QUESTIONS-SCIENCE/2#REPLYTOCOM=116209#RESPOND](https://www.analyticsvidhya.com/blog/2016/09/interview-questions-science/#comment-116209)  
**AVBYTES (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/)**  
**CONTACT US (HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/)**

**WE ARE HIRING! (HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/)**

[\(https://campus.ambassador \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/\)\)](https://campus.ambassador (HTTPS://WWW.ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/))  
 /data-engineering-hackathon-2016/104-hackathon-interview-questions-asked-at-startups-in-machine-6145)

**CORPORATE (HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/)**  
 RECEPTION OF ANSWERS. FROM A  
 NOW THERE ARE SOME  
 IT IS SURELY USEFUL FOR JOB  
 BIGGER FIRMS.

**THE YOUNG  
 OPTIMIZATION  
 CRACKERJACK  
 CONTEST**

By IEOR @ IITB & McKinsey  
 Knowledge Centre India

**REGISTER**

**Win Prizes Worth  
 INR 1 Lakh+**

**content/Team says**  
[INTERVIEW-QUESTIONS-SCIENCE/2#REPLYTOCOM=116150#RESPOND](https://www.analyticsvidhya.com/blog/2016/09/interview-questions-science/#comment-116150)  
**LEARNING- DATA- SCIENCE/ #COMMENT-116150**

[https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/116150/?utm\\_source=AVblog\\_sidebarbottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/116150/?utm_source=AVblog_sidebarbottom) for sharing your thoughts. Tell me more about Q40. What's about it?

**REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING- DATA- SCIENCE/2#REPLYTOCOM=116148#RESPOND)**  
[INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING- DATA- SCIENCE/ #COMMENT-116148](https://www.analyticsvidhya.com/blog/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING- DATA- SCIENCE/ #COMMENT-116148)

I think you got Q3 wrong.

It was to calculate from median and not mean.  
 how can assume mean and median to be same

**Raju (https://www.analyticsvidhya.com/blogs/116149#comment-116149) says:**  
 REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#REPLYTOCOM=116149#RESPOND)  
 AnalyticsVidhya  
 Learn everything about analytics

LEARN ▾ ENGAGE ▾ COMPETE ▾ 

**GET HIRED ▾**  
 Don't bother....Noted ....you assumed normal distribution....



(https://www.analyticsvidhya.com/blogs/116149#comment-116149)  
 /campus-ambassador/

/data-engineering-talent-hunt-hackathon



couraging words! The article is to help beginners tackle the tricky side of ML interviews.

(https://datahack.analyticsvidhya.com/contest/  
 REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#REPLYTOCOM=116161#RESPOND)  
 /the-young-optimization-crackerjack-contest/  
 ?utm\_source=AVBlog\_116161#comment-116161)

Dear Kunal,

Few queries i have regarding AIC

- 1)why we multiply -2 to the AIC equation
- 2)where this equation has been built.

Rgds

**Sampath says:**  
 REPLY (HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#REPLYTOCOM=116195#RESPOND)



/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-

LEARNING-DATA-SCIENCE/#COMMENT-116205

HOME (<https://www.analyticsvidhya.com/>) LEARN ▾ ENGAGE ▾ COMPETE ▾



Hi Manish,

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

**GET HIRED** ▾

Great Job! It is a very good collection of interview

questions. It will be a great help if

**AVBYTES** (<https://www.analyticsvidhya.com/blog/category/abytes/>)

**Data Engineering**

**CONTACT US** (<https://www.analyticsvidhya.com/contact/>)

**Hackathon**

**WE ARE HIRING!** (<https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/>)

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING- DATA-SCIENCE/?REPLYTOCOM=116207#RESPOND)

<https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#comment-116207>

(<https://www.analyticsvidhya.com/campus-ambassador/>)

/data-engineering-talent-hunt-hackathon

Li Comonth



**CORPORATE** (<https://www.analyticsvidhya.com/corporate/>)

Interview questions surely consider  
the following articles.

**THE YOUNG  
OPTIMIZATION  
CRACKER JACK  
CONTEST**

REPLY (<https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/?REPLYTOCOM=116253#RESPOND>)

By IEOR @ IITB & McKinsey  
Knowledge Centre India

**Win Prizes Worth  
INR 1 Lakh+**

**Govkarthik16** says:

Kudos to you all Good Collection for beginners

[https://datahack.kanbay.com/the-young-optimization-crackerjack-contest/?utm\\_source=AnalyticsVidhya](https://datahack.kanbay.com/the-young-optimization-crackerjack-contest/?utm_source=AnalyticsVidhya)

I have small suggestion on Dimensionality Reduction, We can also use the below mentioned techniques to reduce the dimension of the data.

### 1.Missing Values Ratio

Data columns with too many missing values are unlikely to carry much useful information. Thus data columns with number of missing values greater than a given threshold can be removed. The higher the threshold, the more aggressive the reduction.

### 2.Low Variance Filter

Similarly to the previous technique, data columns with little changes in the data carry little information. Thus all data columns with variance lower than a given threshold



are removed. A word of caution: variance is range

dependent therefore normalization is required before

applying this technique.  
(<https://www.analyticsvidhya.com>)

**GET HIRED** ▾

High Correlation Filter



trends are also likely to  
in this case only one of  
machine learning model.

CONTACT US (<https://WWW.ANALYTICSVIDHYA.COM/CONTACT/>)  
**Hackathon**  
WE ARE HIRING! (<https://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/>)

(<https://CAMPUS-AMBASSADOR.ANALYTICSVIDHYA.COM/ANALYTICSVIDHYA.COM/CAMPUS-AMBASSADOR/>)

reduced to anyone. A word of caution: correlation is  
scale sensitive; therefore column normalization is

recommended for a successful correlation comparison.

**CORPORATE** (<https://WWW.ANALYTICSVIDHYA.COM/CORPORATE/>)

Decision Tree, Random Forest, Gradient Boosting Trees

so referred to as random  
selection in addition to  
the approach to  
generate a large and  
trees against a target  
attribute's usage statistics to

([https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm\\_source=AVblog\\_sidebottom](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/?utm_source=AVblog_sidebottom))  
we can generate a large set (2000) of very shallow trees  
(2 levels), with each tree being trained on a small fraction  
(3) of the total number of attributes. If an attribute is often  
selected as best split, it is most likely an informative

feature to retain. A score calculated on the attribute  
usage statistics in the random forest tells us – relative to  
the other attributes – which are the most predictive  
attributes.

## 5. Backward Feature Elimination

In this technique, at a given iteration, the selected  
classification algorithm is trained on n input features.

Then we remove one input feature at a time and train the  
same model on n-1 input features n times. The input  
feature whose removal has produced the smallest  
increase in the error rate is removed, leaving us with n-1



REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-116458](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#COMMENT-116458))  
**Nikhil says:** Analytics Vidhya ([HTTPS://WWW.ANALYTICSVIDHYA.COM/](https://www.analyticsvidhya.com/))

LEARN ▾ ENGAGE ▾ COMPETE ▾ 

### GET HIRED

An awesome article for reference. Thanks a ton Manish sir

pdf format of this blog

**AVBYTES** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/BYTES/](https://www.analyticsvidhya.com/blog/category/abytes/))

**CONTACT US** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/](https://www.analyticsvidhya.com/contact/))

INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-

OND) ([WWW.ANALYTICSVIDHYA.COM](https://www.analyticsvidhya.com))

ASSED-AT-STARTUPS-IN-MACHINE-

15)

**WE ARE HIRING!** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/](https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/))

([HTTPS://CAMPUS.AMBASSADOR.VIDHYA.COM/CAMPUS-AMBASSADOR/](https://campus.ambassador.vidhya.com/campus-ambassador/))  
 /data-engineering-talent-hunt-hackathon  
 also thanks again

**CORPORATE** ([HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/](https://www.analyticsvidhya.com/corporate/))

**THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST**

By IEOR @ IITB & McKinsey  
Knowledge Centre India

**REGISTER**

Win Prizes Worth  
INR 1 Lakh+

0-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-  
OND) ([WWW.ANALYTICSVIDHYA.COM](https://www.analyticsvidhya.com))  
ONS-ASKED-AT-STARTUPS-IN-MACHINE-  
6907)

h. BTW.. I believe the  
ance in question 39 is  
ets are messed. Following

s.

/Bias%E2%80%

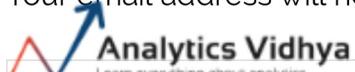
([HTTPS://DATAHACK.ANALYTICSVIDHYA.COM/CONTEST/THE-YOUNG-optimization-crackerjack-contest/](https://datahack.analyticsvidhya.com/contest/the-young-optimization-crackerjack-contest/))  
 /?utm\_source=AVblog\_sidebottom )

REPLY ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/2016/09/40-INTERVIEW-QUESTIONS-ASKED-AT-STARTUPS-IN-MACHINE-LEARNING-DATA-SCIENCE/#COMMENT-118048](https://www.analyticsvidhya.com/blog/2016/09/40-interview-questions-asked-at-startups-in-machine-learning-data-science/#COMMENT-118048))  
**Siddh says:** Analytics Vidhya ([HTTPS://WWW.ANALYTICSVIDHYA.COM/](https://www.analyticsvidhya.com/))

Really awesome article thanks. Given the influence  
young, budding students of machine learning will likely  
have in the future, your article is of great value.

**LEAVE A REPLY**

Your email address will not be published.



[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](#) LEARN ENGAGE COMPETE

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

[GET HIRED](#)

**ABInBev**

[AVBYTES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/AVBYTES/\)](#)

**Data Engineering**

[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT/\)](#)

**Hackathon**

[WE ARE HIRING! \(HTTPS://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/CAREER-ANALYTICS-VIDHYA/\)](#)

(<https://www.analyticsvidhya.com/campus-ambassador>)

/data-engineering-talent-hunt-hackathon

[CORPORATE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CORPORATE/\)](#)

**THE YOUNG  
OPTIMIZATION  
CRACKERJACK  
CONTEST**

By IEOR @ IITB & McKinsey  
Knowledge Centre India

[REGISTER](#)

Win Prizes Worth  
INR 1 Lakh+

[SUBMIT COM](#)

(<https://datahack.analyticsvidhya.com/contest>

/the-young-optimization-crackerjack-contest

?utm\_source=AVblog\_sidebottom )

