

# **DATABRICKS-CERTIFIED- PR OFSSIONAL-DATA-SCIENTIST<sup>Q&As</sup>**

Databricks Certified Professional Data Scientist Exam

**Pass Databricks DATABRICKS-CERTIFIED-  
PROFESSIONAL-DATA-SCIENTIST Exam with 100%  
Guarantee**

Free Download Real Questions & Answers **PDF** and **VCE** file from:

<https://www.leadspass.com/databricks-certified-professional-data-scientist.html>

**100% Passing Guarantee**  
**100% Money Back Assurance**

Following Questions and Answers are all new published by Databricks  
Official Exam Center

- ⚙️ **Instant Download** After Purchase
- ⚙️ **100% Money Back** Guarantee
- ⚙️ **365 Days** Free Update
- ⚙️ **800,000+** Satisfied Customers



**QUESTION 1**

Refer to exhibit

Independent Variable	Coefficient	P-Value
A	0.45	0
B	3.67	0
C	1.23	0

$$R^2 = 0.10$$

You are asked to write a report on how specific variables impact your client's sales using a data set provided to you by the client. The data includes 15 variables that the client views as directly related to sales, and you are restricted to these variables only. After a preliminary analysis of the data, the following findings were made: 1. Multicollinearity is not an issue among the variables 2. Only three variables-A, B, and C-have significant correlation with sales You build a linear regression model on the dependent variable of sales with the independent variables of A, B, and C. The results of the regression are seen in the exhibit. You cannot request additional data. what is a way that you could try to increase the R2 of the model without artificially inflating it?

- A. Create clusters based on the data and use them as model inputs
- B. Force all 15 variables into the model as independent variables
- C. Create interaction variables based only on variables A, B, and C
- D. Break variables A, B, and C into their own univariate models

Correct Answer: A

Explanation: In statistics, linear regression is an approach for modeling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables (or independent variable) denoted  $X$ . The case of one explanatory variable is called simple linear regression. For more than one explanatory variable, the process is called multiple linear regression. (This term should be distinguished from multivariate linear regression<sup>^</sup> where multiple correlated dependent variables are predicted, rather than a single scalar variable.) In linear regression data are modeled using linear predictor functions, and unknown model parameters are estimated from the data. Such models are called linear models. Most commonly, linear regression refers to a model in which the conditional mean of  $y$  given the value of  $X$  is an affine function of  $X$ . Less commonly: linear regression could refer to a model in which the median, or some other quantile of the conditional distribution of  $y$  given  $X$  is expressed as a linear function of  $X$ . Like all forms of regression analysis, linear regression focuses on the conditional probability distribution of  $y$  given  $X$ , rather than on the joint probability distribution of  $y$  and  $X$ : which is the domain of multivariate analysis.

**QUESTION 2**

RMSE is a useful metric for evaluating which types of models?

- A. Logistic regression
- B. Naive Bayes classifier

C. Linear regression

D. All of the above

Correct Answer: C

Explanation: Error calculation allows you to see how well a machine learning method is performing.

One way of determining this performance is to calculate a numerical error. This number is sometimes a percent,

however it can also be a score or distance. The goal is usually to minimize an error percent or distance:

however the goal may be to minimize or maximize a score. Encog supports the following error calculation methods.

Sum of Squares Error (ESS)

Root Mean Square Error (RMS)

Mean Square Error (MSE) (default)

SOM Error (Euclidean Distance Error)

RMSE measures error of a predicted numeric value, and so applies to contexts like regression and some recommender system techniques, which rely on predicting a numeric value. It is not relevant to classification techniques

like logistic regression and Naive Bayes, which predict categorical values. It also is not relevant to unsupervised techniques like clustering. The root-mean-square deviation (RMSD) or root-mean-square error (RMSE) is a frequently used

measure of the

differences between values predicted by a model or an estimator and the values actually observed. Basically,

the RMSD represents the sample standard deviation of the differences between predicted values and observed values.

These individual differences are called residuals when the calculations are performed over the data sample that was used for estimation, and are called prediction errors when computed out-of-sample. The RMSD serves to aggregate the

magnitudes

of the errors in predictions for various times into a single measure of predictive power. RMSD is a good measure of accuracy,

but only to compare forecasting errors of different models for a particular variable and not between variables, as it is scale-dependent.

---

### QUESTION 3

If E1 and E2 are two events, how do you represent the conditional probability given that E2 occurs given that E1 has occurred?

A.  $P(E1)/P(E2)$

B.  $P(E1+E2)/P(E1)$

C.  $P(E2)/P(E1)$

D.  $P(E2)/(P(E1+E2))$

Correct Answer: C

---

## QUESTION 4

In statistics, maximum-likelihood estimation (MLE) is a method of estimating the parameters of a statistical model. When applied to a data set and given a statistical model, maximum-likelihood estimation provides estimates for the model's parameters and the normalizing constant usually ignored in MLEs because:

A. The normalizing constant is always very close to 1

B. The normalizing constant only has a small impact on the maximum likelihood

C. The normalizing constant is often zero and can cause division by zero

D. The normalizing constant doesn't impact the maximizing value

Correct Answer: D

Explanation: (Change the explanation even it is correct)A normalizing constant is positive, and multiplying or dividing a series of values by a positive number does not affect which of them is the largest. Maximum likelihood estimation is concerned only with finding a maximum value, so normalizing constants can be ignored.

---

## QUESTION 5

Suppose you have made a model for the rating system, which rates between 1 to 5 stars. And you calculated that RMSE value is 1.0 then which of the following is correct

A. It means that your predictions are on average one star off of what people really think

B. It means that your predictions are on average two star off of what people really think

C. It means that your predictions are on average three star off of what people really think

D. It means that your predictions are on average four star off of what people really think

Correct Answer: A

[Latest DATABRICKS-CERTIFIED-PROFESSIONAL-DATA-SCIENTIST Dumps](#)

[DATABRICKS-CERTIFIED-PROFESSIONAL-DATA-SCIENTIST VCE Dumps](#)

[DATABRICKS-CERTIFIED-PROFESSIONAL-DATA-SCIENTIST Exam Questions](#)