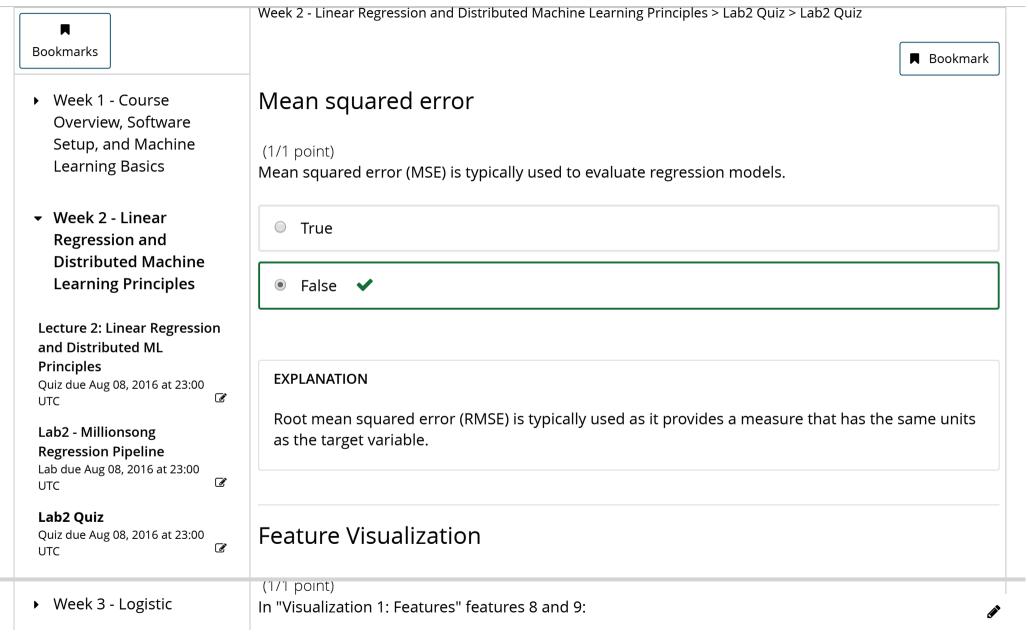


BerkeleyX: CS120x Distributed Machine Learning with Apache Spark



Regression and Clickthrough Rate Prediction

\bigcirc	Have	similar	variation
------------	------	---------	-----------

- Feature 8 varies more than feature 9
- Feature 9 varies more than feature 8

EXPLANATION

Feature 8 is a similar shade of gray throughout while feature 9 varies from light to dark. The standard deviations (across these 50 samples) of feature 8 and feature 9 are .04 and .11, respectively.

Overfitting

(1/1 point)

Select the true statements regarding overfitting:

- Regularization is used to protect against overfitting
- In the lab, we strongly overfit the data



Note: Make sure you select all of the correct options—there may be more than one!

EXPLANATION

Regularization penalizes model complexity, which helps to reduce overfitting. In the lab, our performance on the validation and test sets was comparable so it is unlikely we overfit the data. It appears that our model generalizes well.

Grid Search

(1/1 point)

Select the true statements about grid search:

- lacktriangledown We should conduct a new grid search if we add new features \lacktriangledown
- We can overfit the validation dataset during grid search
- Grid search is computational cheap



Note: Make sure you select all of the correct options—there may be more than one!

EXPLANATION

If the features change, we are fitting a new model, and we should optimize that model through grid search. It is possible that we'll overfit the validation dataset during grid search, which is why we should use a test set to obtain a final model evaluation. Grid search requires training many models, so it is computationally expensive.

Grid Search Visualization

(1/1 point)

In "Visualization 6: Hyperparameter heat map", the best performing models are found in the:

- Top-left quadrant
- Top-right quadrant
- Bottom-left quadrant
- Bottom-right quadrant

EXPLANATION

Model performance was better (as indicated by lighter color) for lower alpha and lower regularization values. This area corresponds to the top-left of the graphs.

Quadratic Features

(1/1 point)

What, if any, impact did quadratic features have on validation error relative to the best model from the grid search you performed?

- None at all
- They increased RMSE by over a year
- They decreased RMSE by over a year

EXPLANATION

RMSE dropped from 15.305 for the grid search result to 14.350, an improvement in RMSE of more than 1 year.

Final Model

(1/1 point)

How did the final model's RMSE compare between the validation and test sets?

RMSE was the same			
RMSE on the validation set was higher			
■ RMSE on the test set was higher			
EXPLANATION			
The RMSE was higher (16.3 vs. 15.7) on the test set compared to the validation set. This increase was similar for the baseline model which had scores of 22.1 and 21.6, respectively.			
Survey: Lab2 Completion Time			
(1/1 point) How long did Lab TWO take you to complete (in hours - decimals are OK)?			
6			
Answer: 0			
$oxed{6}$			
Please click "Check" to save your answer.			

⊚ ③ ⑤ ⑤ Some Rights Reserved



© edX Inc. All rights reserved except where noted. EdX, Open edX and the edX and Open EdX logos are registered trademarks or trademarks of edX Inc.















