



< Previous



Next >

## Programming Assignment 4

🔖 Bookmark this page

Click this link to download the [Diabetes Regression notebook](#) and then complete problems 1-4.

---

Click this link to download the [mystery.dat file](#) which will help you complete problem 5.

---

Click this link to download the [Sentiment Logistic Regression notebook](#).

---

## Problem 1

1/1 point (graded)

This problem is based on the *Diabetes Regression notebook*. You should work through that notebook before entering your answers here.

If a single feature is to be used to predict  $y$ , the best choice (the one that yields the smallest MSE) is feature **2** ('body mass index'). What is the second-best choice? Your answer should be the feature number (**0 – 9**).



Submit

## Problem 2

2/2 points (graded)

Use the **split\_data** procedure to create training/test splits of various sizes. In particular, try training set sizes of **20**, **50**, **100**, and **200**. In each case, record the training error and test error *when using all features for prediction*.

For a training set size of **100**, what are the training MSE and test MSE (just round to the nearest integer)?

Training MSE =



Test MSE =



Submit

## Problem 3

1/1 point (graded)

What *rough* trends do you observe as the training set size increases (from, say, **20** to **400**)? Select all that apply.

☒ The training error increases

☒ The test error decreases

☒ The gap between the training and test error decreases



Submit

## Problem 4

1/1 point (graded)

What is the single best explanation for these trends? Choose one of the following.

☐ With more training data, we get better estimates of training error.

☒ With more training data, we learn a more accurate model.

☐ The error is proportional to the amount of data.



Submit

Problem 5 relates to finding relevant features.

## Problem 5

1/1 point (graded)

The file **mystery.dat** contains pairs  $(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{x} \in \mathbb{R}^{100}$  and  $\mathbf{y} \in \mathbb{R}$ . There is one data point per line, with comma-separated values; the very last number in each line is the  $\mathbf{y}$ -value.

In this data set,  $\mathbf{y}$  is a linear function of just *ten* of the features in  $\mathbf{x}$ , plus some noise. Your job is to identify those ten features.

Which of the following contain only relevant features?

(Think of the feature numbers as being in the range 1 to 100, but be aware that Python indexes arrays starting at zero.)

☐ 1,5,7,19,44

☒ 2,3,13,17,29

☐ 3,7,13,19,44

☐ 5,23,24,51,61



Submit

Submit

< Previous

Next >

© All Rights Reserved



## edX

[About](#)

[Affiliates](#)

[edX for Business](#)

[Open edX](#)

[Careers](#)

[News](#)

## Legal

[Terms of Service & Honor Code](#)

[Privacy Policy](#)

[Accessibility Policy](#)

[Trademark Policy](#)

[Sitemap](#)

## Connect

[Blog](#)

[Contact Us](#)

[Help Center](#)

[Media Kit](#)



© 2022 edX LLC. All rights reserved.

深圳市恒宇博科技有限公司 [粤ICP备17044299号-2](#)