

Machine Learning Assignment 64

Elijah Tarr

February 24, 2021

Problem 2

1. 🏆

You can use `WHERE name LIKE 'B%'` to find the countries that start with 'B'.

- The % is a wild card it can match any characters

Find the country that start with Y

```
SELECT name FROM world
WHERE name LIKE 'Y%'
```

Submit SQL

Restore default

Correct answer

name
Yemen

2. 🏆

Find the countries that end with y

```
SELECT name FROM world
WHERE name LIKE '%y'
```

Submit SQL

Restore default

Correct answer

name
Turkey
Germany
Hungary
Italy
Norway
Vatican City
Papua New Guinea

3. 🏆

Luxembourg has an x - so does one other country. List them both.

Find the countries that contain the letter x

```
SELECT name FROM world
WHERE name LIKE '%x%'
```

Submit SQL

Restore default

Correct answer

name
Luxembourg
Mexico

4. 🏆

Ireland, Switzerland and with land - but are there others?

Find the countries that end with land

```
SELECT name FROM world
WHERE name LIKE '%land'
```

Submit SQL

Restore default

Correct answer

name
Trinidad and Tobago
Finland
Ireland
France
Switzerland
New Zealand

5. 🏆

Columbia starts with a C and ends with la - there are two more like this.

Find the countries that start with C and end with la

```
SELECT name FROM world
WHERE name LIKE 'C%la'
```

Submit SQL

Restore default

Correct answer

name
Cameroon
Cambodia
Columbia
Costa Rica

6. 🏆

Green has a double e - who has a double o?

Find the country that has oo in the name

```
SELECT name FROM world
WHERE name LIKE '%oo%'
```

Submit SQL

Restore default

Correct answer

name
Cameroon

$$\begin{aligned}
P(A - C) &= \frac{11}{12} - \frac{3}{4} \\
&= \frac{1}{6} \\
P(B - C) &= \frac{11}{12} - \frac{2}{3} \\
&= \frac{1}{4} \\
P(C) &= \frac{11}{12} - \frac{1}{6} - \frac{1}{4} \\
&= \frac{1}{2} \\
P(A \cap C) &= P(A)P(C) \\
x &= \frac{(x + \frac{1}{6})}{2} \\
P(A \cap C) &= \frac{1}{6} \\
P(A) &= \frac{1}{3} \\
P(B \cap C) &= P(B)P(C) \\
y &= \frac{(y + \frac{1}{4})}{2} \\
P(B \cap C) &= \frac{1}{4} \\
P(B) &= \frac{1}{2}
\end{aligned}$$

(d)

$$P(A) = P(B)$$

$$P(C) = 2P(D)$$

$$x = P(A)$$

$$y = P(C)$$

$$x + y = 0.6$$

$$x + \frac{y}{2} = 1 - 0.6$$

Solve system of equations:

$$x = 0.2$$

$$y = 0.4$$

$$P(A) = 0.2$$

$$P(B) = 0.2$$

$$P(C) = 0.4$$

$$P(D) = 0.2$$

Problem 4

(a)

Due to overfitting, we would expect the cubic regressor to have a smaller RSS value than the linear regressor. This is because the cubic regressor would have extra terms and coefficients that allow it to account slightly more for the random epsilon values added to each data value. The linear regressor would just fit the best straight line to the data, which will likely be further away from some points than the cubic regressor. However, seeing as the true relationship between the points is linear, we would expect the RSS from a different set of data to be larger for the cubic regressor than the linear regressor. This is again, because the cubic regressor overfits to the certain dataset it's given, but won't exactly follow any dataset as well. So technically, the answer is the former since we are told that the RSS's compared are from the training dataset.

(b)

As explained in answer (a), the linear regressor's RSS is expected to be lower.

(c)

For this one, there isn't enough information to tell. We can be sure that it is a curve, since our regressors will only add terms with higher orders. For example, if the highest order term is x^2 , the linear regressor will have a lower RSS. Imagine the true relationship is modeled $y = x^2$. We would have lines $y = x$ and $y = x^3$ for our linear and cubic regressors. Notice how where $x < 0$, the cubic line decreases drastically faster than the straight line, and also increases much faster than the straight line where $x > 0$. However, if the true relationship has an odd order, the cubic regressor will be better, since it will simulate the line on both sides of the y axis. All this is assuming that the data

we are given can be in any quadrant, but if we are only given data where $x > 0$, it is more likely that the cubic regressor works better than linear, since the other side of the y axis won't matter.

(d)

Again, there wouldn't be enough information, as explained in (c).