Computational Statistics

(1) This is a preview of the published version of the quiz

Started: Oct 19 at 9:08am

Quiz Instructions

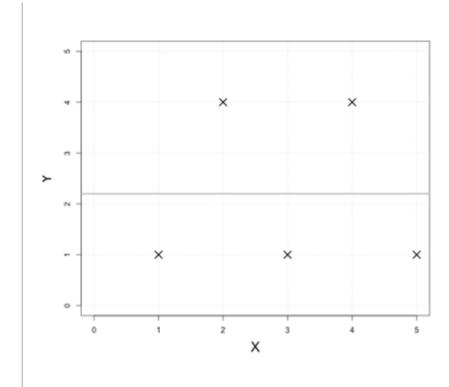
Aids:

- Non-programmable calculator

Hint:

- Use period as decimal separator, not comma

Data frame *df1* contains 5 points: $(x_1,y_1)=(1,1), (x_2,y_2)=(2,4), (x_3,y_3)=(3,1), (x_4,y_4)=(4,4),$ and $(x_5,y_5)=(5,1)$ as shown here:





What is the MSE (mean squared error) of the model shown as the gray line at $\hat{y}=2.2$ for the entire data set? Round your answer to two decimal places.

2.16





Consider df1 again. You randomly sampled points 2, 3, and 5 into the validation set. What is the validation MSE of the linear regression model predicting Y from X? Provide your calculations for partial credit.

Edit View Insert Format Tools Table

Training set contains (1,1) and (4,4). Linear fit will be y = x.

Validation MSE = $((2-4)^{**2} + (3-1)^{**2} + (5-1)^{**2})/3 = (4+4+16)/3 = 8$

р







You collect a set of data points (n = 100 observations) containing a single feature X and a quantitative target Y. You first split the data into 50% training data *train* and 50% validation data *validation*. Using the training data and Python, you then fit two regression models: a **linear regression model** using only the single, non-transformed feature X, and a **cubic regression model** using the original feature X as well as additional, transformed features X² and X³.

Question 3



Suppose that the **true relationship between X and Y is linear**, i.e., $Y=b_1X+b_0+\epsilon$ for some b_0,b_1 and with some Gaussian noise ϵ with $\mathbf{E}[\epsilon]=\mathbf{0}$. Consider the **training MSE** for the linear regression and the training MSE for the cubic regression. Would we expect:

- There is not enough information to tell which regression model has lower error.
- O The errors of both regression models to be about the same.
- The error of the cubic regression to be lower.
- The error of the linear regression to be lower.

Question 4

1.5 pts

Suppose that the true relationship between X and Y is still linear but consider now the **validation MSE** for the linear regression and the cubic regression. Would we expect:

- O There is not enough information to tell which regression model should have lower error.
- O The error of the cubic regression to be lower.
- O The errors of both regression models to be about the same.
- The error of the linear regression to be lower.

Question 5



Suppose now that the **true relationship between X and Y is not linear**, **but we don't know how far it is from linear**. Consider the **training MSE** for the linear regression and the cubic regression. Would we expect:

- O The errors of both regression models should be about the same.
- The error of the cubic regression to be lower.
- \bigcirc The error of the linear regression to be lower.
- There is not enough information to tell which regression model should have lower error.

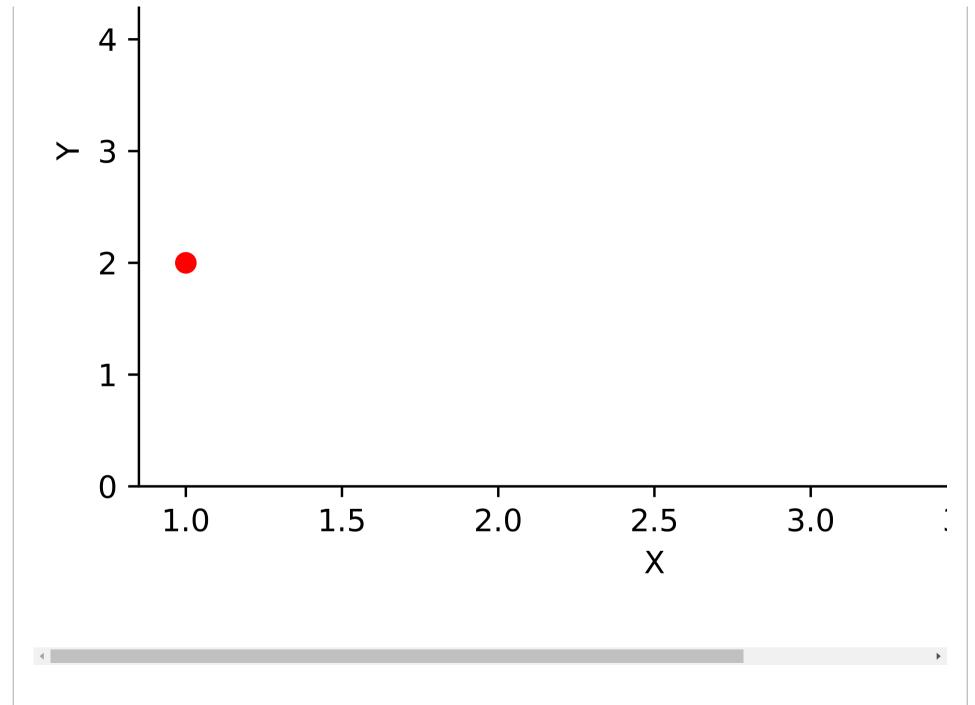
1.5 pts

Suppose the true relationship between X and Y is still not linear and we still don't know how far it is from linear. Consider the **validation MSE** for the linear regression and the cubic regression. Would we expect:

- O The error of the cubic regression to be lower.
- There is not enough information to tell which regression model should have lower error.
- O The errors of both regression models should be about the same.
- The error of the linear regression to be lower.

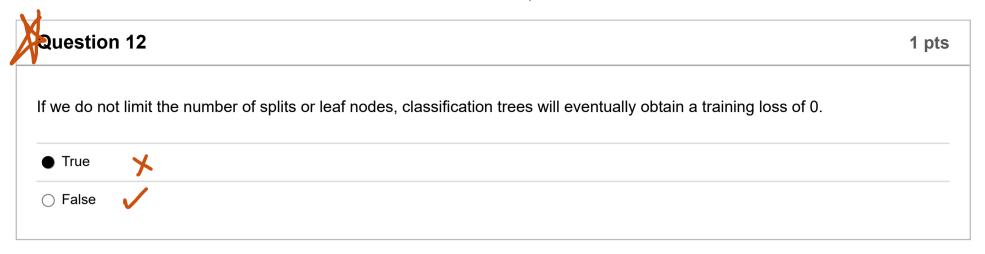
When growing decision trees, the greedy recursive binary splitting algorithm needs to find the optimal split with respect to a given loss function. Consider the following example of a regression task with a single feature X and a target Y:

5 -



Question 7	Type text here	3 pts
What is the SSE (sum of squared e	errors) on this training data for the split at X=1.5?	
Question 8		3 pts
What is the SSE on this training da	ta for the split at X=3.0?	
Question 9		2 pts
What is the SSE when splitting at b	ooth X=1.5 and X=3.0 on this training data?	

Multiple Choice Questions:
Question 10
The maximal margin classifier guarantees the lowest test error among all possible hyperplanes when the training data is linearly separable.
○ True● False
Question 11 1 pts
In a classification setting, if the data is not linearly separable, one should extend the feature space by adding transformed features. Exchanging the original features with transformed features cannot lead to a better fit to the data.
○ True
● False





When growing a random forest, for each tree, one first samples a subset of features, and then restricts each split in that tree to only using one of those tree-specific subsampled features (with replacement).

○ True

False

Question 14



Random forests with base models that are grown by searching over every feature for each split (max_features=n_features in sklearn.ensemble.RandomForestClassifier) are just bagged trees.

True			
○ False			

Random forests determine feature importance by looking at how often that feature appears at the root node. True False

You have a dataset cancer_data.csv that contains data on patients who have gotten a sample of breast mass taken with each row representing one patient. There are 4 features and 1 target variable in the columns. The features describe characteristics of the cell nuclei present in the image of those samples and the target denotes whether a patient's sample is malignant (0) or benign (1). The objective is to train a random forest classifier using the features to predict the probability that the sample is benign. To evaluate your model, you are using the Brier (quadratic) score and the validation set method.

Large parts of the code are already given. Your task is to fill in the missing pieces. The parts with missing pieces are marked in yellow (the length of the yellow parts are not necessarily the same as the length of the missing code):

1 import pandas as pd

```
IMPORT NUMBY as no
     from sklearn.model_selection import train_test_split
     from sklearn.ensemble import RandomForestClassifier
     def brier(
                                ):
         return 1-((y-model.
                                        (X)[
                                                       1)**2)
 8
     df = pd.read_csv("cancer_data.csv", index_col=0)
10
11
                          [:,4]
                 = df.
                 = df. [:.:4]
12
13
14
     X_train, X_validate, y_train, y_validate = train_test_split(features, targets, test_s
15
16
17
18
19
     print("Average Brier score of Random Forest on validation set:", brier(rf,
```



Please enter (only!) the missing code for line 6:

model, X, y



Please enter (only!) the missing code for the first yellow part of line 7:

predict_proba

Question 18



Please enter (only!) the missing code for the **second yellow part of line 7** (you can assume that the model was trained on a dataset containing both classes):

:,-1

Question 19



Please enter (only!) the missing code for the **third yellow part of line 7**:

.mean()

Question 20	1 pts
Please enter (only!) the missing code for the first yellow part of line 11 :	
targets	
Question 21	0.5 pts
Please enter (only!) the missing code for the second yellow part of line 11 :	
iloc	
Question 22	1 pts
Please enter (only!) the missing code for the first yellow part of line 12:	
features	

0.5 pts

Please enter (only!) the missing code for the **second yellow part of line 12**:

iloc

Question 24



Please enter the missing code for line 16 (creating the random forest object; you can leave all parameters at the default values):

rf = RandomForestClassifier()

Question 25



Please enter the missing code for line 17 (training the model):

rf.fit(X_train, y_train)

1 pts

Please enter (only!) the missing code for **line 19**:

X validate, y validate

You are a probabilistic forecaster competing with two other forecasters for a single prize that will be given to the forecaster with highest quadratic scoring rule score. (Ties will be broken uniformly at random, so that all forecasters with highest quadratic score receive the prize with equal probability.)

The quadratic scoring rule is: $R_{q}\left(y,\,x
ight)\,=\,1-\left(y-x
ight)^{2}$

The forecasting competition has only a single question and you believe that the event will happen with probability p = 0.6. You know that the two other forecasters also believe that the probability of the event occurring is 0.6. Moreover, you know that they will report truthfully.

Question 27

1 pts

What is your subjective probability of being selected if you report 0.6 truthfully?

○ 2/3			
○ None of these			
<u> </u>			

Question 28	2 pt
Which of these forecasts maximizes your subjective probability of winning the prize?	
O.6	
● 1.0	
○ None of these	
O Any of these	
O.0	

2 pts

Which of these forecasts maximizes your subjective probability of winning the prize if the other two forecasters are not truthful but report 1.0 (i.e., 100%) instead?

○ None of these

O.6

Any of these



0.0



Question 30



What is the highest expected probability of winning the prize that you can obtain assuming the other two forecasters report 1.0 (i.e., 100%)?

0.4

Question 31

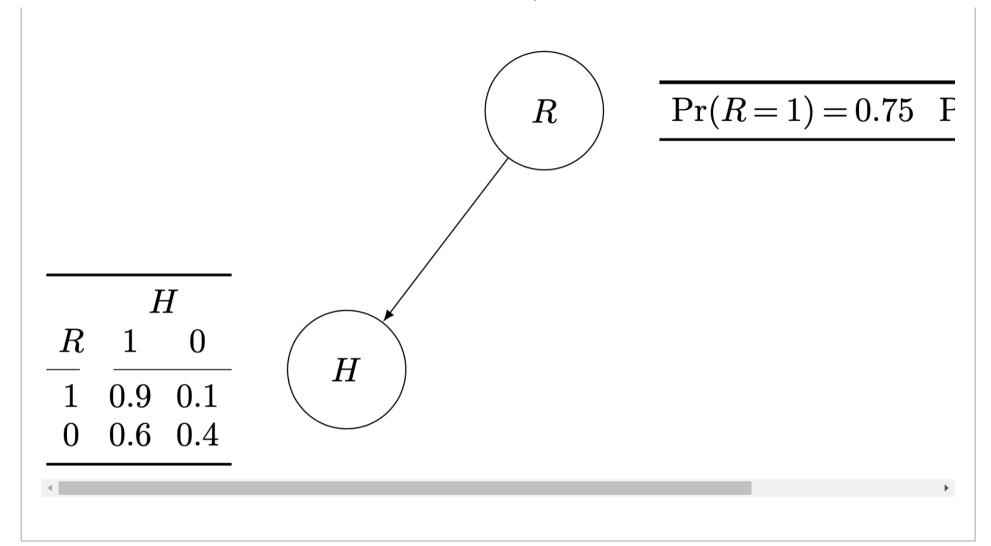


In the lecture, we have seen that we can implement truthful forecasting competitions by giving the prize to forecaster i with probability

$$f_i = rac{1}{n} + rac{1}{n} \Biggl(R_q(y_i,x) - rac{1}{n-1} \sum_{j
eq i} R_q(y_j,x) \Biggr) \,.$$

For x=1, compute the probability that forecaster 1 receives the prize when she reported $y_1=0.5$ and the other two forecasters reported $y_2=y_3=1.0$.

Consider the following Bayesian network representing the beliefs of a weather forecaster seeking to predict the probability of rain (R). The forecaster has been offered access to a (noisy) hygrometer (H) with the conditional probability table shown in the graphic.





What is the forecaster's belief that the hygrometer signal will be 1? That is, what is Pr(H=1)? (use **at least three decimals** in your calculations wherever possible)

0.825

Question 33



What is the forecaster's belief of rain given that the hygrometer signal is 0? That is, what is Pr(R=1|H=0)? (use **at least three decimals** in your calculations wherever possible)

0.429

Question 34



What is the expected quadratic score of the truthfully-reporting forecaster if she does not learn the outcome of H (i.e., if she does not have access to the hygrometer)? (The numerical result alone is sufficient but you can provide your calculation for partial credit; use **at least three decimals** in your calculations wherever possible.)

Edit View Insert Format Tools Table

 $I_{\mathcal{D}} \quad \blacksquare \quad \vee \quad \sqrt{\mathsf{x}} \quad \Leftrightarrow$

$$E = 0.75*(1-(1-0.75)**2 + 0.25*(1-0.75**2)$$

 $E = 0.8125$

р





0 words | </> /





Question 35



The forecaster still seeks to predict the probability of R=1 and is still reporting truthfully. Her expected quadratic score following H=1 is 103/121 = 0.851, her expected quadratic score following H=0 is 37/49 = 0.755.

What is the forecaster's expected improvement in quadratic score from learning the outcome of H? (The numerical result alone is sufficient but you can provide your calculation for partial credit; use at least three decimals in your calculations wherever possible.)

Edit View Insert Format Tools Table

0.825*103/121 + 0.175*37/49 = 0.8344

Expected improvement of 0.8344-0.8125 = 0.0219





0 words | </> </





Not saved

Submit Quiz