DSC 80 Discussion 8 Worksheet

READ THIS if you're completing this worksheet to submit to Gradescope for discussion credit:

- You don't need to fill in your answers directly on this worksheet, but all 7 questions are numbered, so if you're writing your answers elsewhere, make sure it's clear what's answering what.
- Either way, your responses should include at least a couple sentences of explanation for every question in order to receive credit.
- Submit your work to Gradescope by **Tuesday**, **May 28**, **by 11:59 PM**. There will be no extensions or late submissions accepted for this assignment. (It's optional, anyhow.)
- If you're having trouble, you can watch the Friday discussion podcast, although we might not get to every question during that time. If you really aren't sure, try your best as long as you legitimately attempt every question, you'll get credit or describe what you're having trouble with.

1 FA22 Final Problem 7

	group	color	x	у
0	Α	red	3	2
1	В	green	7	1
2	Α	blue	2	5
3	Α	red	5	3
4	В	blue	10	4
5	Α	green	1	1

Consider the dataframe to the left. Suppose you wish to use this data in a linear regression model. To do so, the color column must be encoded numerically.

Problem 1.1. True or False: a meaningful way to numerically encode the color column is to replace each string by its index in the alphabetic ordering of the colors. That is, to replace blue by 1, green by 2, and red by 3.

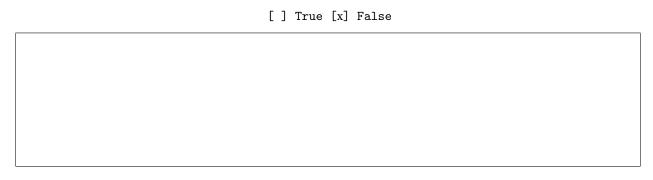
[]	True
[x]	False

Problem 1.2. scikit-learn's OneHotEncoder module has a keyword called drop-first, which the documentation says will "drop the first category in each feature." What's the purpose of this keyword, and will using it lead to a worse linear classifier?

Reduce redundancy, avoid multicollinearity

2 FA22 Final Problem 9

Problem 2.1. Suppose you split a data set into a training set and a test set. You train a classifier on the training set and test it on the test set. **True or False**: the training accuracy must be higher than the test accuracy.



Problem 2.2. Suppose you train a model, but achieve much lower training and test accuracies than you expect. When you look at the data and make predictions yourself, you are easily able to achieve higher train and test accuracies. What should be done to improve the performance of the model?

Note: You haven't learned about decision trees yet (basically, just imagine a flow-chart), but for this question, all you need to know is that increasing max_depth increases the complexity of your model.

[]	Decrease the max_depth hyperparameter; the model is "over	erfitting".
$\lceil x \rceil$	Increase the max_depth hyperparameter: the model is "und	erfitting".

3 SP22 Final Problem 10

The DataFrame new_releases contains the following information for songs that were recently released. The first few rows are shown below.

	genre	rec_label	danceability	speechiness	first_month
0	Hip-Hop/Rap	ЕМІ	0.39	0.84	12019896
1	Pop	UMG	0.91	0.65	9932385
2	Pop	EMI	0.65	0.71	10923584
3	Country	SME	0.45	0.93	8107742
4	Hip-Hop/Rap	UMG	0.39	0.86	9554136

- genre: one of the following five possibilities: Hip-Hop/Rap, Pop, Country, Alternative, or International
- rec_label: the label that released the song (one of the following 4: EMI, SME, UMG, or WMG)

- danceability: how easy the song is to dance to, according to the Spotify API (between 0 and 1)
- speechiness: what proportion of the song is made up of spoken words, according to the Spotify API
- first_month: the number of total streams the song had on Spotify in the first month it was released

To start, we conduct a train-test split, splitting new_releases into X_train, X_test, y_train, and y_test. We first fit a linear model to the training data that only uses danceability, and call this model lr_one.

Problem 3.1. True or False: If lr_one.score(X_train, y_train) is much lower than lr_one.score(X_test, y_test), it is likely that lr_one overfit to the training data.

```
[ ] True [x] False
```

```
>>> X_train.shape[0]
50
>>> np.sum((y_train - lr_one.predict(X_train)) ** 2)
500000 # five hundred thousand
```

Problem 3.2. Given this output, what is lr_one's training RMSE? Give your answer as an integer.

```
100
```

Now, suppose we fit one more linear model (with an intercept term) to the training data:

• Model 2 (lr_no_drop): Uses danceability and speechiness as-is, and one-hot encodes genre and rec_label, using OneHotEncoder(). (Note the lack of the drop_first=True keyword.)

Suppose we are given the following coefficients in Model 2:

- The coefficient on genre_Pop is 2000.
- The coefficient on genre_Country is 1000.
- The coefficient on danceability is $10^6 = 1,000,000$

Problem 3.3. Daisy and Billy are two artists signed to the same rec_label who each just released a new song with the same speechiness. Daisy is a Pop artist while Billy is a Country artist.

Model 2 predicted that Daisy's song and Billy's song will have the same first_month streams. What is the absolute difference between Daisy's song's danceability and Billy's song's danceability? Give your

answer as a simplified fi	raction.		
1/1000			

DSC 80 Discussion 9 Worksheet

1 WI24 Final Problem 10

Harshi is trying to build a decision tree that predicts whether or not a hotel has a swimming pool. Suppose + represents the "has pool" class and \circ represents the "no pool" class. One of the nodes in Harshi's decision tree has 12 points, with the following distribution of classes.

Consider the following three splits of the node above.

- Split 1: Yes: $++\circ\circ\circ$ No: $++\circ\circ\circ$
- Split 2: Yes: $++\circ\circ\circ\circ\circ\circ$ No: $++\circ$
- Split 3: Yes: \circ No: $++++\circ\circ\circ\circ\circ\circ$

Which of the six nodes above have the same entropy as the original node? **Select all that apply.**

Which of the six nodes above have the lowest entropy? (There can be multiple correct answers.)

- [x] Split 1's "Yes" node
- [x] Split 1's "No" node
- [] Split 2's "Yes" node
- [x] Split 2's "No" node
- [] Split 3's "Yes" node
- [] Split 3's "No" node

- [] Split 1's "Yes" node
- [] Split 1's "No" node
- [] Split 2's "Yes" node
- [] Split 2's "No" node
- [x] Split 3's "Yes" node
- [] Split 3's "No" node

2 WI24 Final Problem 12

Diego wants to build a model that predicts the number of open rooms a hotel has, given various other features. He has a training set with 1200 rows. Diego fits a regression model using the GPTRegression class.

GPTRegression models have several hyperparameters that can be tuned, including context_length and sentience. To choose between 5 possible values of context_length, Diego performs k-fold cross-validation.

How many total times is a GPTRegression model fit?

[] 4k

[x] 5k

[] 240k

[] 6000k

[] 4(k-1)

[] 5(k-1)

[] 240(k-1)

[] 6000(k-1)

Suppose that every time a GPTRegression model is fit, it appends the number of points in its training set to the list sizes. Note that after performing cross-validation, len(sizes) is equal to your answer to the previous subpart.

What is sum(sizes)?

[] 4k

[] 5k

[] 240k

[] 6000k

[] 4(k-1)

[] 5(k-1)

[] 240(k-1)

[x] 6000(k-1)

3 SP22 Final Problem 12

After fitting a classifier, we use it to make predictions on an unseen test set. Our results are summarized in the following confusion matrix.

Actually Negative ??? 30 Actually Positive 66 105

What is the recall of our classifier? Give your answer as a fraction (it does not need to be simplified).
105/171
The accuracy of our classifier is $\frac{69}{117}$. How many true negatives did our classifier have?
33
True or False: In order for a binary classifier's precision and recall to be equal, the number of mistakes makes must be an even number.
Тrue

Suppose we are building a classifier that listens to an audio source (say, from your phone's microphone) and predicts whether or not it is Soulja Boy's 2008 classic "Kiss Me thru the Phone." Our classifier is pretty good at detecting when the input stream is "Kiss Me thru the Phone", but it often incorrectly predicts that similar sounding songs are also "Kiss Me thru the Phone."

Our classifier has:

[] low precision and low recall.

[x] low precision and high recall.

[] high precision and low recall.

[] high precision and high recall.

4 WI23 Final Problem 8.4

Let's say we're performing k-fold cross-validation to find the best combination of two hyperparameters called height and color.

For the purposes of this question, assume that:

- Our training set contains n rows, where n is greater than 5 and is a multiple of k.
- There are h_1 possible values of height and h_2 possible values of color.

Solve for the following quantities, in terms of n, k, h_1 , and h_2 :

- What is the size of each fold? n/k
- How many times is row 5 in the training set used for training? __h1h2(k-1)
- How many times is row 5 in the training set used for validation? h1h2