

DS-100 Final Exam

Spring 2017

Name: _____

Email address: _____

Student id: _____

Instructions:

- Please fill in your name, email address, and student id at the top of both this exam booklet and your answer sheet.
- All answers must be written on the separate answer sheet.
- This exam must be completed in the **3 hour time** period ending at **6:00PM**.
- You may use a single page (two-sided) cheat sheet.
- Work quickly through each question. There are a total of 175 points on this exam.
- You must turn in both this exam booklet and your answer sheet.
- **Don't cheat!**

1 Maximum Likelihood and Loss Minimization

1. Suppose we observe a dataset $\{x_1, \dots, x_n\}$ of independent and identically distributed samples from the exponential distribution. The probability density function (PDF) of an exponential distribution (for $x \geq 0$) parameterized by the parameter λ is given by:

$$f_\lambda(x) = \lambda e^{-\lambda x}$$

- (1) [4 Pts.] What is the *log-likelihood function* of this *dataset* with respect to λ ?

Solution: The likelihood function is given by:

$$L(\lambda) = \prod_{i=1}^n f_\lambda(x_i) = \prod_{i=1}^n \lambda e^{-\lambda x_i} \quad (1)$$

Therefore the log-likelihood function is given by:

$$\log L(\lambda) = \sum_{i=1}^n \log(\lambda e^{-\lambda x_i}) = \sum_{i=1}^n \log(\lambda) + \log(e^{-\lambda x_i}) \quad (2)$$

$$= n \log(\lambda) - \lambda \sum_{i=1}^n x_i \quad (3)$$

- (2) [6 Pts.] Derive the maximum likelihood value $\hat{\lambda}_{\text{MLE}}$. **Circle your answer.**

Solution: Taking the derivative of the log-likelihood function with respect to the parameter λ we get:

$$\frac{\partial}{\partial \lambda} \log L(\lambda) = n \frac{\partial}{\partial \lambda} \log(\lambda) - \frac{\partial}{\partial \lambda} \lambda \sum_{i=1}^n x_i \quad (4)$$

$$= n \frac{1}{\lambda} - \sum_{i=1}^n x_i \quad (5)$$

$$(6)$$

To compute the maximum likelihood parameter $\hat{\lambda}_{\text{MLE}}$ we set the above derivative equal

to zero and solve.

$$0 = n \frac{1}{\hat{\lambda}_{\text{MLE}}} - \sum_{i=1}^n x_i \quad (7)$$

$$\frac{1}{\hat{\lambda}_{\text{MLE}}} = \frac{1}{n} \sum_{i=1}^n x_i \quad (8)$$

$$\hat{\lambda}_{\text{MLE}} = \frac{n}{\sum_{i=1}^n x_i} \quad (9)$$

Thus the maximum likelihood parameter estimate is:

$$\hat{\lambda}_{\text{MLE}} = \frac{n}{\sum_{i=1}^n x_i} = \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^{-1} = \frac{1}{\mathbf{Mean}(x)} \quad (10)$$

You may use the following scratch space but we will only grade what you put on the answer sheet.

2. Suppose we collect a dataset of n IID observations $\{x_1, \dots, x_n\}$ which we believe are drawn from a distribution with the following PDF:

$$f_\mu(x) = C \exp\left(-\frac{(x - \mu)^6}{6}\right) \quad (11)$$

where C is a constant that does not depend on μ .

- (1) [3 Pts.] Write the log-likelihood function for μ .

Solution: Because we assumed the data are independent and identically distributed (IID) according to f_μ , the likelihood function is the product of the probabilities of each observation:

$$L(\mu) = \prod_{i=1}^n f_\mu(x_i) = C^n \prod_{i=1}^n \exp\left(-\frac{(x_i - \mu)^6}{6}\right) \quad (12)$$

$$= C^n \exp\left(-\frac{1}{6} \sum_{i=1}^n (x_i - \mu)^6\right) \quad (13)$$

Taking the log of the likelihood function we get:

$$\log L(\mu) = n \log C - \frac{1}{6} \sum_{i=1}^n (x_i - \mu)^6 \quad (14)$$

- (2) [4 Pts.] Compute the derivative of the log-likelihood with respect to μ .

Solution: Taking the derivative:

$$\frac{\partial}{\partial \mu} \log L(\mu) = \frac{\partial}{\partial \mu} n \log C - \frac{1}{6} \sum_{i=1}^n \frac{\partial}{\partial \mu} (x_i - \mu)^6 \quad (15)$$

$$= 0 + \sum_{i=1}^n (x_i - \mu)^5 \quad (16)$$

$$= \sum_{i=1}^n (x_i - \mu)^5 \quad (17)$$

- (3) [3 Pts.] Because there is no closed form solution for μ in $\frac{\partial}{\partial \mu} \log L(\mu) = 0$, we would likely use gradient *ascent* to approximately compute $\hat{\mu}_{\text{MLE}}$. Given the gradient function:

$$g(\mu) = \frac{\partial}{\partial \mu} \log L(\mu), \quad (18)$$

and a step size $\rho(t)$, what is the gradient ascent update rule to go from $\mu^{(t)}$ to $\mu^{(t+1)}$? (*Hint: your answer should contain only the variables $g(\mu^{(t)})$, $\mu^{(t)}$, $\mu^{(t+1)}$, and $\rho(t)$.*)

Solution: Recall that the gradient points in the “*uphill*” direction. When we maximize, uphill is the way to want to go. So the update rule would look like:

$$\mu^{(t+1)} \leftarrow \mu^{(t)} + \rho(t)g(\mu^{(t)}) \quad (19)$$

2 Wrangling and Querying Data

2.1 SQL

For the questions in this subsection, assume we have a massive database in the cloud with the following schema:

```
-- A simple digital media store database
CREATE TABLE media
  (mid integer PRIMARY KEY,
   name text, type char, year_released integer, length integer,
   buy_cost float, rent_cost float, avg_rating float);

CREATE TABLE customers
  (cid integer PRIMARY KEY,
   name text, joined date, nation_id integer,
   activity_level integer);

CREATE TABLE transactions
  (tid integer PRIMARY KEY,
   tdate date, item integer, customer integer,
   rent_or_buy integer, price_paid float, percent_viewed float,
   FOREIGN KEY (item) REFERENCES media,
   FOREIGN KEY (customer) REFERENCES customers);

CREATE VIEW stats AS
SELECT min(length) AS len_min, max(length) AS len_max,
       avg(length) AS len_mu, stddev(length) AS len_sigma,
       min(avg_rating) AS ar_min, max(avg_rating) AS ar_max,
       avg(avg_rating) AS ar_mu, stddev(avg_rating) AS ar_sigma
FROM media;
```

3. [4 Pts.] In the `media` table above, the `type` column encodes the type of media as a unique character code (e.g., 'S' for song, 'M' for movie, 'E' for episode, etc.). Suppose we wanted to modify the `stats` view to display the stats for each `type` of media. Which of the following are true? (Select *all* that apply.)

- A. We need to change the granularity of the view to be finer than it is above.
- B. We need to add a `GROUP BY type` clause to the view.
- C. It would be helpful to add `media.type` to the list of columns in the `SELECT` clause of the view.
- D. The modified view should have more rows than the original view above.
- E. None of the above.

4. [3 Pts.] Which of the following queries finds the ids of media that are at least 2 standard deviations longer than the mean length? (Select *only one*.)

A.

```
SELECT media.mid
FROM media, stats
WHERE media.mid = stats.mid
      AND media.length >= stats.len_mu
                          + 2*(stats.len_sigma);
```

B.

```
SELECT media.mid
FROM media, stats
WHERE media.length >= stats.len_mu
                          + 2*(stats.len_sigma);
```

C.

```
SELECT media.mid
FROM media
WHERE media.length >= avg(media.length)
                          + 2*stddev(media.length);
```

D. None of the above.

2.2 SQL Sampling

The `transactions` table has 30 million (30×10^6) rows. It is too large to load into the memory of our laptop. We will extract a sample from the database server to process on our laptop in Python.

```
SELECT *  
FROM transactions TABLESAMPLE Bernoulli(.0001);
```

5. [2 Pts.] Suppose you ran this query many times. What distribution describes the output sizes (in number of rows) you would see across runs?

Solution: Binomial

6. [2 Pts.] In expectation, how many rows will there be in the answer to this query?

Solution: $3000 = 30 \times 10^2$

7. [4 Pts.] Your friend Emily Engineer tells you to avoid Bernoulli sampling, and use the following query instead:

```
SELECT *  
FROM transactions  
LIMIT XX;
```

(where `XX` is replaced by the correct answer to the previous question). **Select all the true statements:**

- A. Emily's `LIMIT` query will probably run faster than the `TABLESAMPLE` query. For Emily's query, the database engine can simply access the first `XX` rows it finds in the table, and skip the rest.**

Solution: True. For reasoning above.

- B. Emily's query result may be biased to favor certain rows.**

Solution: True. The database will optimize for speed, which will likely favor clusters of records stored near each other.

- C. The output of the `TABLESAMPLE` query provides a hint about how many rows there are in the `transactions` table while Emily's `LIMIT` query does not.**

Solution: True. You can extrapolate from the sample size and the sample probability to predict the table size.

- D. Emily's `LIMIT` query may run fast, but it will swamp the memory on your laptop, since it doesn't sample the database.

Solution: False. Emily's query will only return XX rows to the laptop.

- E. None of the above.
8. [2 Pts.] You will recall from Homework 4 that it is possible to do bootstrap sampling in SQL by constructing a `design` table with two columns. Each of the columns used in that scheme is described by a single choice below. **Identify the *two* correct choices:**
- A. A foreign key to the table being sampled.**
 - B. A `count` column to capture the number of tuples in each bootstrap sample.
 - C. An identifier to group rows together into bootstrap samples.**
 - D. A regularization column to prevent overfitting.

2.3 Pandas

For the questions in this subsection, assume that we have pandas dataframes with the same schemas as described in the previous section on SQL. That is, we have a `media` dataframe with columns `mid`, `name`, `type`, `year`, et cetera. Assume that the index column of each dataframe is meaningless—the primary key is represented as a regular column.

9. [3 Pts.] Consider the following code snippet:

```
def get_average_price_paid(join_method):
    return (customers
            .merge(transactions, how=join_method,
                  left_on='cid', right_on='customer')
            .loc[:, 'price_paid']
            .fillna(0)           # <- Important
            .mean()
    )

inner = get_average_price_paid('inner')
outer = get_average_price_paid('outer')
left = get_average_price_paid('left')
right = get_average_price_paid('right')
```

Assume that all item *prices are positive*, all `transactions` refer to valid customers in the `customers` table, but some customers may have no transactions.

- (1) How are `inner` and `outer` related? **Pick one best answer.**
- A. `inner < outer`
 - B. `inner ≤ outer`
 - C. `inner = outer`
 - D. `inner ≥ outer`**
 - E. `inner > outer`
- (2) How are `left` and `right` related? **Pick one best answer.**
- A. `left < right`
 - B. `left ≤ right`**
 - C. `left = right`
 - D. `left ≥ right`
 - E. `left > right`
- (3) How are `left` and `outer` related? **Pick one best answer.**
- A. `left < outer`
 - B. `left ≤ outer`
 - C. `left = outer`**
 - D. `left ≥ outer`
10. [3 Pts.] We wish to write a python expression to find the largest amount of money spent by one person on any single date. We will use the following code:
- ```
biggie = transactions.groupby(____)['price_paid'].sum().max()
```
- What should we be pass in as our `groupby` predicate? **Select only one answer.**
- A. `'tdate'`
  - B. `'customer'`
  - C. `['item', 'tdate']`
  - D. `['customer', 'tdate']`**
  - E. `['customer', 'item']`
11. [6 Pts.] Fill in the following python code that finds the names of every customer who has spent over \$100.
- ```
merged = customers.merge(__A__, left_on=__B__, right_on=__C__)
grouped = merged.groupby(__D__).__E__()
names = grouped[__F__].index
```

Solution:

```
merged = customers.merge(transactions, \
                          left_on="cid", right_on="customer")
grouped = merged.groupby("cid").sum()
names = grouped[grouped.price_paid > 100].index
```

12. [4 Pts.] Find the earliest year where the average price_paid that year exceeds the average price_paid over all years. We have the following code:

```
merged = transactions.merge(media, left_on="item", \
                             right_on="mid")
mean_price = merged.groupby("year_released") \
                  .mean().price_paid.mean() # Line A
by_year = merged.groupby("year_released").count() # Line B
is_greater = by_year[by_year.price_paid > mean_price] # Line C
result = is_greater.sort_index(ascending=False).index[0] # Line D
```

Some of these lines need to be modified in order for the code to work properly. We have suggested replacements for each line below. Which lines need to be *replaced*? **Select all that apply.**

A. mean_price = merged.price_paid.mean()

B. by_year = merged.groupby("year_released").mean()

C. is_greater = by_year.where(by_year.price_paid > mean_price)

D. result = is_greater.sort_index(ascending=True).index[0]

E. All the lines are correct.

3 Feature Engineering

For this problem we collected the following data on the new social networking app *UFace*.

PostID	UTC Time	Text	Num. Responses	State
3	08:10 PM	"Checkout my breakfast ..."	2	VA
13	11:00 AM	"Studied all night for ..."	5	CA
14	12:04 PM	"Hello world!"	0	NY
17	11:35 PM	"That exam was lit ..."	42	CA
...				

13. Suppose we are interested in predicting the number of responses *for future posts*. For each of the columns, indicate which (**one or more**) of the given feature transformations could be informative. **Select *all* that apply.**

(1) [2 Pts.] The `PostID` column:

- A. Drop the column**
- B. One-Hot encoding
- C. Leave as is

(2) [2 Pts.] The `Time` column:

- A. Take the hour as a float (between 0 and 24)**
- B. One-Hot encoding
- C. Bag-of-words encoding
- D. Time since midnight in seconds**

(3) [2 Pts.] The `Text` column:

- A. The length of the text**
- B. One-Hot encoding
- C. Bag-of-words encoding**
- D. Leave as is

(4) [2 Pts.] The `State` column:

- A. The length of the text
- B. One-Hot encoding**
- C. Bag-of-words encoding
- D. Leave as is

14. [4 Pts.] Suppose we believe that people are more likely to respond to tweets in the *afternoon* (roughly from hours 13 to 17). Which of the following feature functions would help capture this intuition? Assume that the function **localHour** takes a time and a state as its arguments and returns the hour of the day (in 24-hour time) in the state's time zone. Also assume that any boolean-valued feature is encoded as 0 (false) or 1 (true). **Select all that apply.**
- A. $\phi(\text{time}, \text{state}) = \text{localHour}(\text{time}, \text{state})$
 - B. $\phi(\text{time}, \text{state}) = 13 < \text{localHour}(\text{time}, \text{state}) < 17$
 - C. $\phi(\text{time}, \text{state}) = \exp(-(\text{localHour}(\text{time}, \text{state}) - 15)^2)$
 - D. $\phi(\text{time}, \text{state}) = \exp(\text{localHour}(\text{time}, \text{state}) - 15)$
 - E. None of the above.
15. [2 Pts.] Given the following text from a BigData Borat post:

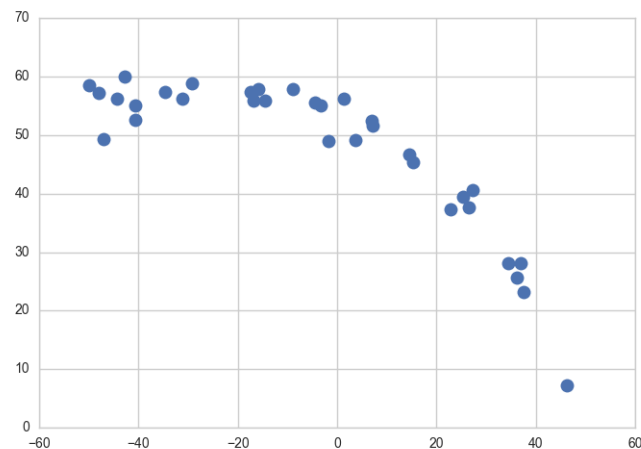
“Data Science is statistics on a Mac.”

Which of the following is the *bi-gram* encoding *including stop-words*? (Select *only one*.)

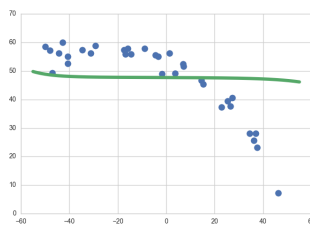
- A. $\{('data', 1), ('science', 1), ('statistics', 1), ('mac', 1)\}$
- B. $\{('data science', 1), ('science statistics', 1), ('statistics mac', 1)\}$
- C. $\{('data science', 1), ('science is', 1), ('is statistics', 1), ('statistics on', 1), ('on a', 1), ('a mac', 1)\}$
- D. $\{('data science', 1), ('is statistics', 1), ('on a', 1), ('mac', 1)\}$

4 Least Squares Regression and Regularization

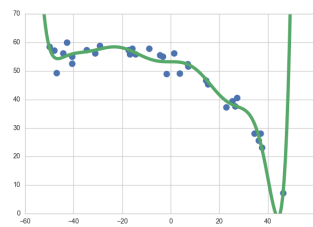
16. For this question we use the following toy dataset:



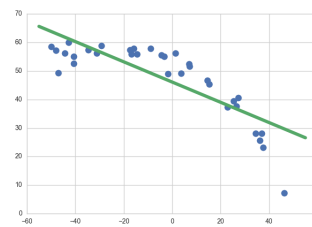
(1) [3 Pts.] We have fit several models depicted as curves in the following plots:



(a)



(b)



(c)

Select the plot that best matches each of the models below. **Each plot is used exactly once.**

1. Linear regression model

☐ (A) ☐ (B) ☒ (C)

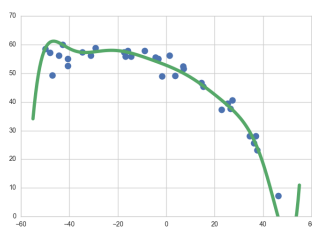
2. Linear regression with degree 10 polynomial features

☐ (A) ☒ (B) ☐ (C)

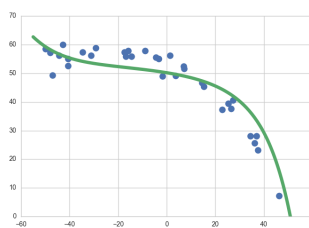
3. Ridge regression with degree 10 polynomial features and substantial regularization.

☒ (A) ☐ (B) ☐ (C)

- (2) [2 Pts.] We fit two more models to these data. Again, the solid curves display the predictions made by each model.



(a)



(b)

Select the plot that best matches each of the models below. **Each plot is used exactly once.**

1. Ridge regression with degree 10 polynomial features, $\lambda = 0.1$.
☒ (A) ☐ (B)
 2. Ridge regression with degree 10 polynomial features, $\lambda = 1.0$.
☐ (A) ☒ (B)
17. Suppose you are given a dataset $\{(x_i, y_i)\}_{i=1}^n$ where $x_i \in \mathbb{R}$ is a one dimensional feature and $y_i \in \mathbb{R}$ is a real-valued response. To model this data you choose a model characterized by the following objective function:

$$J(\theta) = \sum_{i=1}^n (y_i - \theta_0 - x_i\theta_1 - x_i^2\theta_2)^2 + \lambda \sum_{i=1}^2 |\theta_i| \quad (20)$$

- (1) [7 Pts.] **Select all the true statements** for the above objective function (Equation 20).

- A. This loss function likely corresponds to a classification problem.
- B. θ is the regularization parameter.

C. This is an example of L_1 regularization.

- D. This is not a linear model in θ .

E. This model includes a bias/intercept term.

F. This model incorporates a non-linear feature transformation.

G. Large values of λ would reduce the model to a constant θ_0 .

- H. None of the above are true.

- (2) [2 Pts.] Suppose in our implementation we accidentally forget to square the first term:

$$J(\theta) = \sum_{i=1}^n (y_i - \theta_0 - x_i\theta_1 - x_i^2\theta_2) + \lambda \sum_{i=1}^2 |\theta_i| \quad (21)$$

What would change if we tried to train a model using gradient descent on this objective function rather than the original objective function? **(Select only one)**

- A. The training code would raise an error due to a matrix/vector dimension problem.
 - B. The training process would diverge with $\theta_0 \rightarrow -\infty$
 - C. The training process would diverge with $\theta_0 \rightarrow \infty$**
 - D. The training process would converge to a different regression line.
 - E. Nothing; the training process would eventually converge to the same regression line.
18. [5 Pts.] Let X be a $n \times p$ design matrix with full column rank and y be a $n \times 1$ response vector. Let $\hat{\beta}$ be the optimal solution to the least squares problem and r be its associated error. In other words,

$$y = X\hat{\beta} + r \quad (22)$$

Consider X_2 the second column of X .

- (1) [1 Pt.] **True or False.** Without any additional assumptions,

$$r \cdot X_2 = 0$$

where \cdot denotes the usual dot product?

- (2) [4 Pts.] Provide a short proof or counter example.

You may use the following scratch space but we will only grade what you put on the answer sheet.

Solution: True. It suffices to show that r is orthogonal to the column space of X .

$$X^T r = X^T (y - X(X^T X)^{-1} X^T y) = (X^T - X^T X (X^T X)^{-1} X^T) y = (X^T - X^T) y = 0$$

5 Classification

19. For each of the following select **T** for true or **F** for false on the answer sheet.

- (1) [1 Pt.] A binary or multi-class **classification** technique should be used whenever there are **categorical features**.

Solution: False. Categorical *features* may appear in both classification and regression settings and should be addressed using one-hot-encoding.

- (2) [1 Pt.] Logistic regression is actually used for classification.

Solution: True. Logistic regression is somewhat confusingly named as it applies to classifications tasks but builds on the linear models we introduced in least squares linear regression.

- (3) [1 Pt.] The logistic regression loss function was derived by modeling the observations as noisy observations with a Gaussian noise model.

Solution: False. Logistic regression was derived using the Bernoulli likelihood of function.

- (4) [1 Pt.] Class imbalance can be a serious problem in which the number of training data points from one class is much larger than another.

Solution: True. Class imbalance can be a serious problem and often occurs in settings like disease diagnosis where a large fraction of the population is healthy.

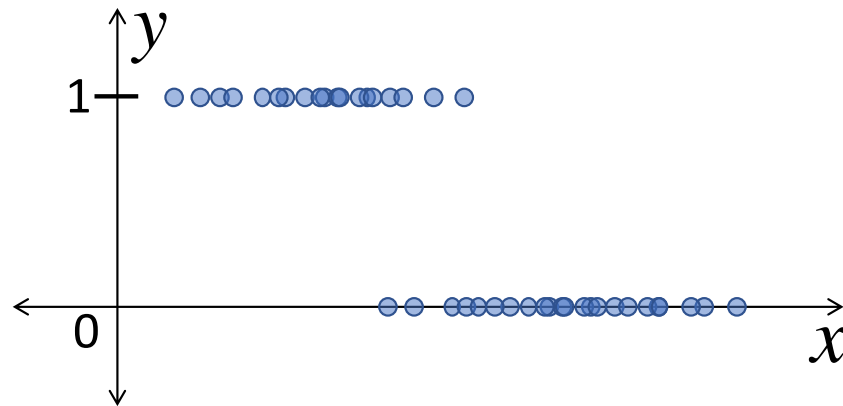
- (5) [1 Pt.] A broken *binary* classifier that *always* predicts 0 is likely to get a test accuracy around 50% on all prediction tasks.

Solution: False. In many case class imbalance could result in substantially higher or lower accuracy.

- (6) [1 Pt.] The root mean squared error is the correct metric for evaluating the prediction accuracy of a binary classifier.

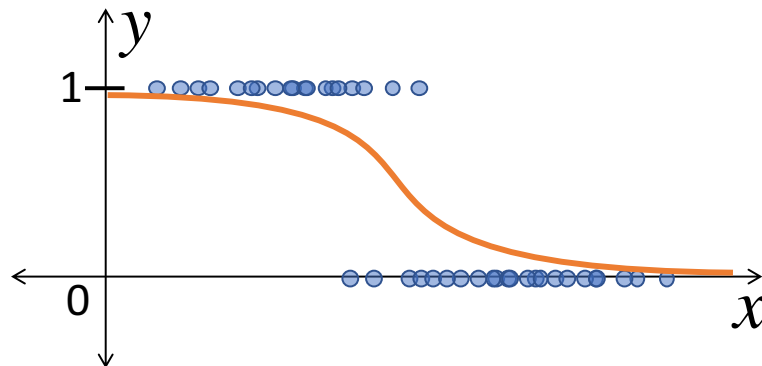
Solution: False. Root mean squared error is a standard measure of accuracy for regression. Logistic regression accuracy is often measured by the fraction of examples predicted correctly or in some cases the likelihood of the data under the model.

20. Consider the following binary classification dataset



- (1) [3 Pts.] Draw a reasonable approximation of the logistic regression probability estimates for $\mathbf{P}(Y = 1 \mid x)$ on top of the figure on the answersheet.

Solution: Anything close to the following would be acceptable:



It is important that:

1. the curve is higher for smaller values of x
2. the curve is smooth
3. the curve is a sigmoid

- (2) [1 Pt.] Are these data linearly separable?

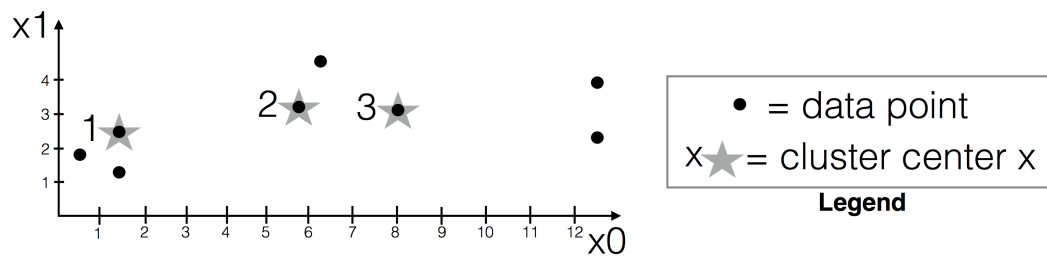
A. Yes

B. No

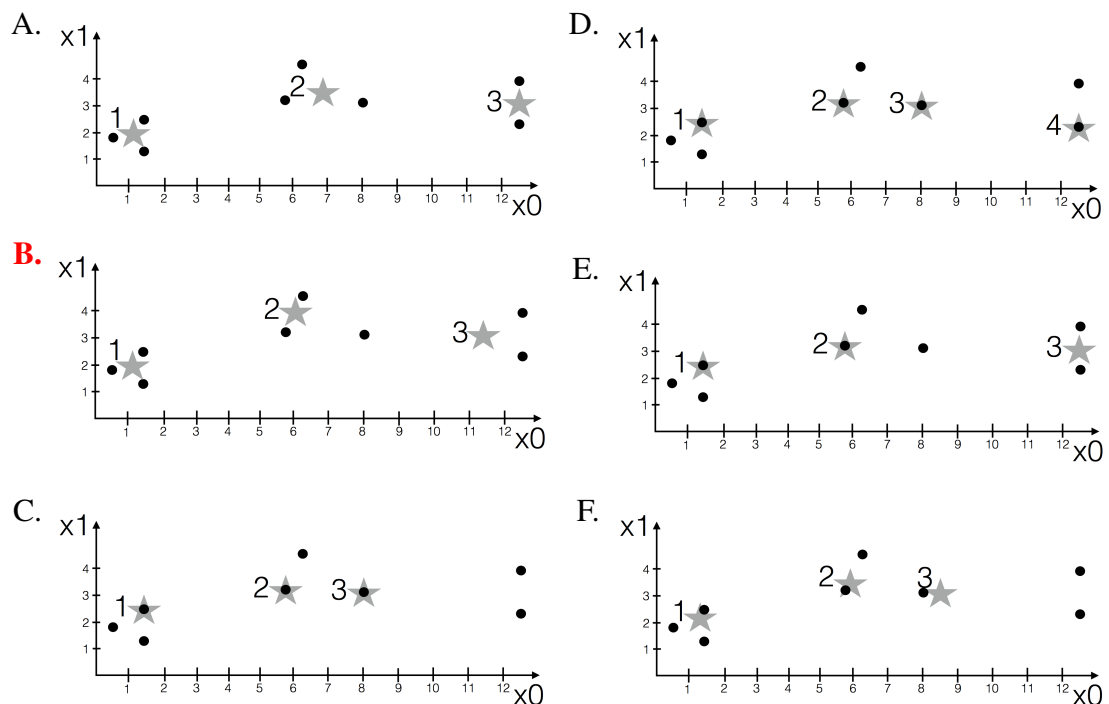
21. [3 Pts.] Suppose you are given θ for the logistic regression model to predict whether a tumor is malignant ($y = 1$) or benign ($y = 0$) based on features of the tumor x . If you get a new patient x_* and find that $x_*^T \theta > 0$, what can you say about the tumor? **Select only one.**
- A. The tumor is benign
 - B. The tumor is more likely benign
 - C. The tumor is more likely to be malignant**
 - D. The tumor is malignant
22. [4 Pts.] Which of the following explanations justify applying regularization to a logistic regression model? **Select all that apply.**
- A. The training error is too high.
 - B. The test error is too low.
 - C. The data are high-dimensional.**
 - D. There is a large class imbalance.
 - E. None of the above justify regularization for logistic regression.

6 Clustering

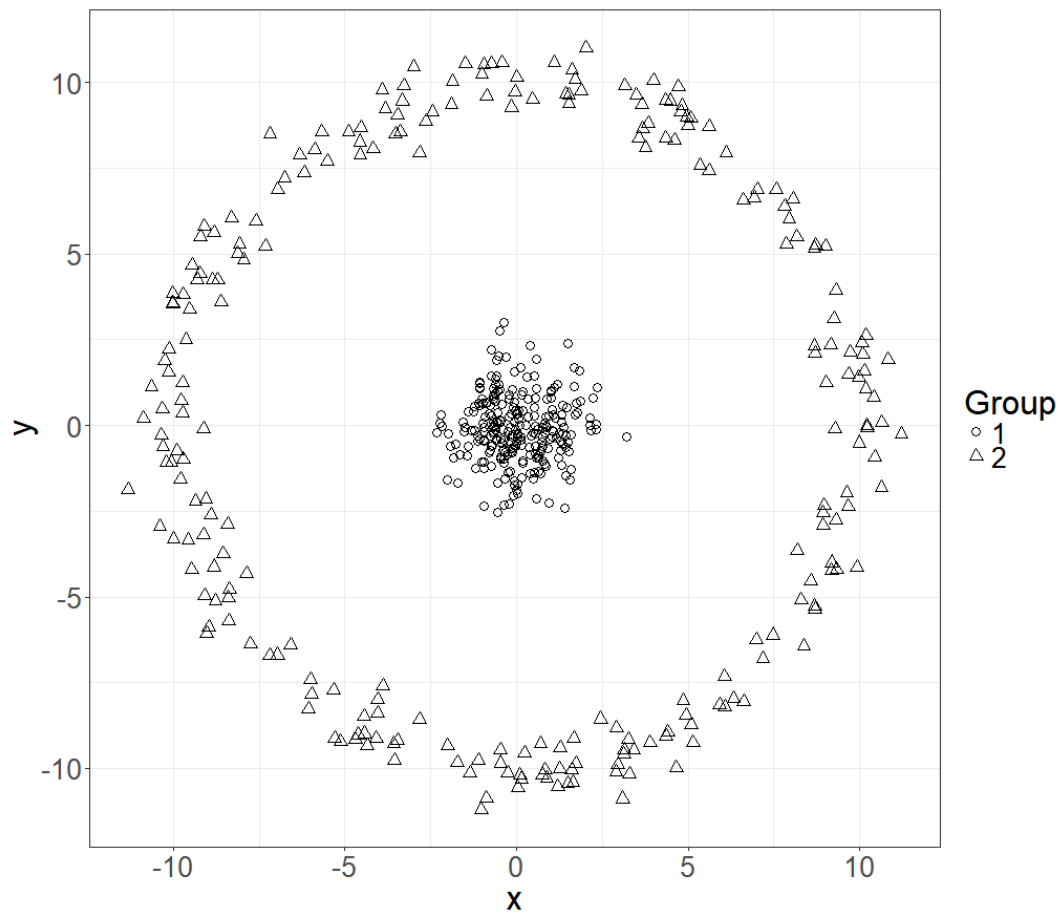
23. [4 Pts.] The following diagram shows a scatter plot of a small 2-dimensional dataset with 8 elements. An *initialization* of the k-means algorithm (with $k = 3$) is displayed; the initial cluster centers are displayed as stars. (They have the same locations as 3 of the points themselves, which is a common initialization of the k-means algorithm.)



Which of the following depicts the 3 cluster centers that would result from *a single iteration* of the k-means algorithm (with $k = 3$), starting from the initial cluster center locations above? **Select only one.**



24. [4 Pts.] Consider the data plotted below. Which of the following clustering methods is most likely to return the groupings below using untransformed x and y values? **Select *only one*.**



- A. Single-linkage clustering
- B. Complete-linkage clustering
- C. Average-linkage clustering
- D. k-means clustering
- E. None of the above are likely to recover the true groups

7 Bias-Variance Tradeoff

25. For each of the following select **T** for true or **F** for false on the answer sheet.

- (1) [1 Pt.] Regularization can be used to manage the bias-variance trade-off.

Solution: True. Regularization encourages simpler models which can help to reduce variance but increase bias.

- (2) [1 Pt.] When conducting linear regression, adding polynomial features to your data often decreases the variance of your fitted model.

Solution: False. Adding more features tends to increase your model's variance since there are more parameters to fit.

- (3) [1 Pt.] When conducting linear regression, adding polynomial features to your data often decreases the bias of your fitted model.

Solution: True. Adding more features tends to decrease your model's bias since your model can fit more complicated patterns in the data.

- (4) [1 Pt.] Suppose your data are an i.i.d. sample from a population. Then collecting a larger sample for use as a *training set* can help reduce *bias*.

Solution: False. Increasing the dataset size without changing the modeling procedure can often reduce variance but is unlikely to address bias.

- (5) [1 Pt.] Suppose your data are an i.i.d. sample from a population. Then collecting a larger sample for use as a *training set* can help reduce *variance*.

Solution: True. More data often helps to reduce variance in the model fitting process.

- (6) [1 Pt.] Training error is typically larger than test error.

Solution: False. Training error often under-estimates the test error.

- (7) [1 Pt.] If you include the test set in your training data, your accuracy as measured on the test set will probably increase.

Solution: True. Training on the test data improves test accuracy but this improvement can be misleading due to over-fitting.

- (8) [1 Pt.] It is important to frequently evaluate models on the test data throughout the process of model development.

Solution: False. Nooooooooooooo. Once test data is used it is no longer test data. You should create validation datasets or use cross-validation procedures to evaluate models.

26. [2 Pts.] A colleague has been developing models all quarter and noticed recently that her *test* error has started to gradually increase while her training error *has been decreasing*. Which of the following is the most likely explanation for what is happening? **Select *only one*.**

- A. She is starting to over-fit to her training data.**
- B. She is starting to under-fit to her training data.
- C. The model is overly biased.
- D. None of the above.

27. [5 Pts.] Given the following general loss formulation:

$$\arg \min_{\theta} \left[\sum_{i=1}^n (y_i - x_i^T \theta)^2 + \lambda \sum_{p=1}^d \theta_p^2 \right] \quad (23)$$

Which of the following statements are true? **Select *all* that apply.**

- A. There are d data points.
- B. There are n data points.**
- C. The data is d dimensional.**
- D. This is a classification problem.
- E. This is a linear model.**
- F. This problem has LASSO regularization.
- G. Larger values of λ imply increased regularization.**
- H. Larger values of λ will increase variance.
- I. Larger values of λ will likely increase bias.**
- J. None of the above are true.

28. [3 Pts.] In class we broke the least-squares error into three separate terms:

$$\mathbf{E} [(y - f_{\theta}(x))^2] = \mathbf{E} [(y - h(x))^2] + \mathbf{E} [(h(x) - f_{\theta}(x))^2] + \mathbf{E} [(f_{\theta}(x) - \mathbf{E}[f_{\theta}(x)])^2] \quad (24)$$

where $y = h(x) + \epsilon$, $h(x)$ is the true model and ϵ is zero-mean noise. For each of the following terms, indicate its usual interpretation in the bias variance trade-off:

1. $\mathbf{E} [(y - h(x))^2]$: A. Bias B. Variance **C. Noise**
2. $\mathbf{E} [(h(x) - f_{\theta}(x))^2]$: **A. Bias** B. Variance C. Noise
3. $\mathbf{E} [(f_{\theta}(x) - \mathbf{E}[f_{\theta}(x)])^2]$: A. Bias **B. Variance** C. Noise

8 Big Data

29. Consider the following simple Data Warehouse schema from a Cellular Service Provider, which records activity on a cell phone network:

```
CREATE TABLE devices (  
    did integer, customer_id integer,  
    phone_number varchar(13),  
    firstname text, lastname text,  
    zip varchar(12), registered_on varchar(2),  
    PRIMARY KEY (did),  
    UNIQUE (customer_id) -- a ``candidate`` key  
);
```

```
CREATE TABLE billing (  
    rate_code char PRIMARY KEY,  
    description text, base_fee float, per_minute float,  
    max_minutes integer, overage_fee float,  
    PRIMARY KEY (rate_code));
```

```
CREATE TABLE calls (  
    caller_handset_id integer, callee_handset_id integer,  
    cell_tower_id integer, call_start datetime, call_end datetime,  
    billing_code char,  
    PRIMARY KEY (caller_handset_id, call_start),  
    FOREIGN KEY (caller_handset_id) REFERENCES devices,  
    FOREIGN KEY (billing_code) REFERENCES billing;
```

- (1) [3 Pts.] Which of these tables is a dimension table? **Select *all* that apply.**
- A. devices
 - B. calls
 - C. billing
 - D. None of the above.

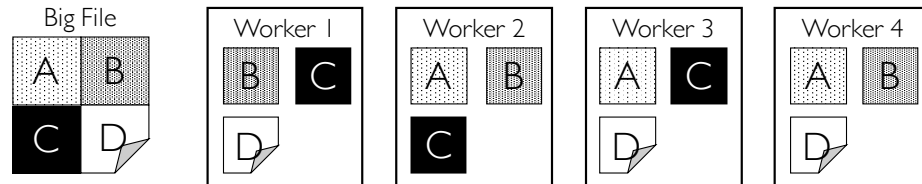
Solution: devices and billing

- (2) [3 Pts.] Which of the following statements are true? **Select *all* that apply.**
- A. The `calls.billing_code` column violates star schema design because any update to a single billing fee requires updates to many call records.
 - B. If we want to look for correlations between a device's average call length and the time since it was registered, we have to perform a join.**

C. If the cell service provider implemented a Data Lake, it would make it easier for them to load audio recordings of calls for subsequent analysis.

D. None of the above statements are true.

30. [3 Pts.] The figure below depicts a distributed file system with one logical “big file” partitioned into 4 “shards” (A, B, C, D) and replicated across multiple worker machines (1, 2, 3, 4).



Suppose workers 1 AND 2 both fail. Which of the following statements are true? **Select all that apply.**

A. The full file will remain available since worker 3 and worker 4 are both still running.

B. The system can tolerate one more worker failure without losing data.

C. If every request requires all 4 shards of the file, then worker 3 and worker 4 can share the work evenly.

D. None of the above statements are true.

31. [3 Pts.] Consider only the mechanism of *partitioning* files into shards, and storing different shards on different machines. Which of the following statements are true? **Select all that apply.**

A. Partitioning enhances the ability of the system to store large files.

B. Partitioning allows the system to tolerate machine failures without losing data.

C. Partitioning allows the system to read files in parallel.

D. None of the above statements are true.

32. [2 Pts.] Recall the statistical query pattern discussed in class for computing on very large data sets. Which of the following statements are true? **Select all that apply.**

A. It eliminates the need for the end-user device (e.g. a laptop) to acquire all the data.

B. It pushes the computational task closer to the large-scale data storage.

C. It is well suited to both MapReduce and SQL interfaces.

D. An alternative to the statistical query pattern for big data is to acquire a sample of the full dataset on the end-user device.

E. None of the above statements are true.

9 EDA and Visualization

33. [2 Pts.] Consider the following statistics for infant mortality rate. According to these statistics, which transformation would best symmetrize the distribution? (**Select only one.**)

Transformation	lower quartile	median	upper quartile
x	13	30	68
\sqrt{x}	3.5	5	8
$\log(x)$	1.15	1.5	1.8

- A. no transformation
 B. square root
C. log
 D. not possible to tell with this information
34. [5 Pts.] For each of the following scenarios, determine which plot type is *most* appropriate to reveal the distribution of and/or the relationships between the following variable(s). **For each scenario, select only one plot type. Some plot types may be used multiple times.**

- | | |
|--------------------------|------------------------|
| A. histogram | F. scatter plot |
| B. pie chart | G. stacked bar plot |
| C. bar plot | H. overlaid line plots |
| D. line plot | I. mosaic plot |
| E. side-by-side boxplots | |

- (1) [1 Pt.] sale price and number of bedrooms (assume integer) for houses sold in Berkeley in 2010.

Solution: E. Side-by-side Boxplots. We might imagine using a scatter plot since we are plotting the relationship between two numeric quantities. However because the number of bedrooms is an integer and most houses will only have a small number, we are likely to encounter *over-plotting* in the scatter plot. Therefore side-by-side boxplots are likely to be most informative.

- (2) [1 Pt.] sale price and date of sale for houses sold in Berkeley between 1995 and 2015.

Solution: F. Scatter Plot. Here we are plotting two numeric quantities with sufficient spread on each axis.

- (3) [1 Pt.] infant birth weight (grams) for babies born at Alta Bates hospital in 2016.

Solution: A. Histogram. Here we are plotting the distribution of a likely large number of observations and therefore a histogram would be most appropriate.

- (4) [1 Pt.] mother's education-level (highest degree held) for students admitted to UC Berkeley in 2016

Solution: C. Bar Plot. Here we want to visualize counts of a categorical variable.

- (5) [1 Pt.] SAT score and HS GPA of students admitted to UC Berkeley in 2016

Solution: F. Scatter Plot. Here we are visualizing the relationship between two continuous quantities.

- (6) [1 Pt.] race and gender of students admitted to UC Berkeley in 2016

Solution: I. mosaic plot Here we are visualizing the relationship between two categorical variables.

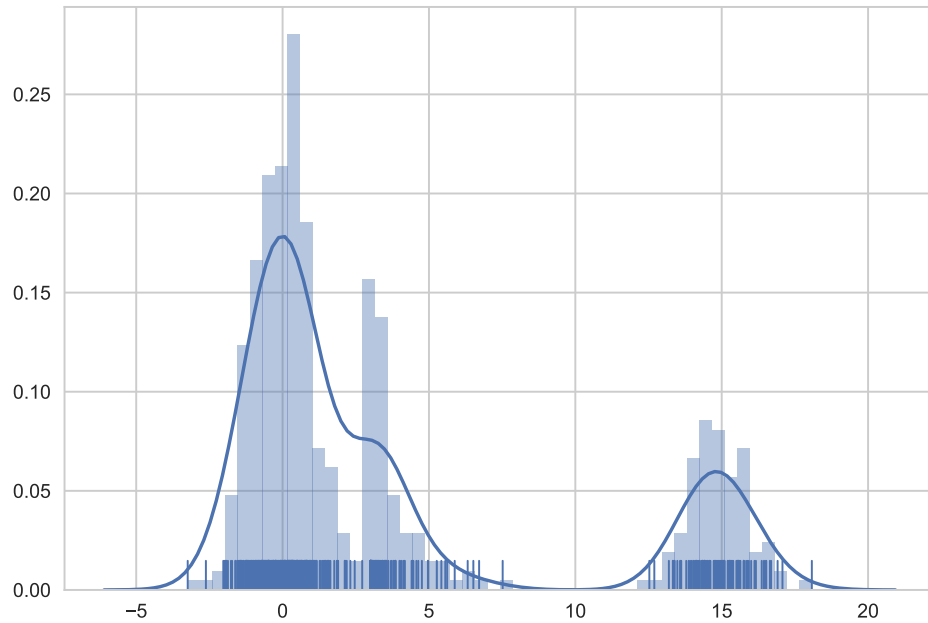
- (7) [1 Pt.] The percentage of female student admitted to UC Berkeley each year from 1950 to 2000.

Solution: D. Line plot. This allows us to see the trends over time.

- (8) [1 Pt.] SAT score for males and females of students admitted to UCB from 1950 to 2000

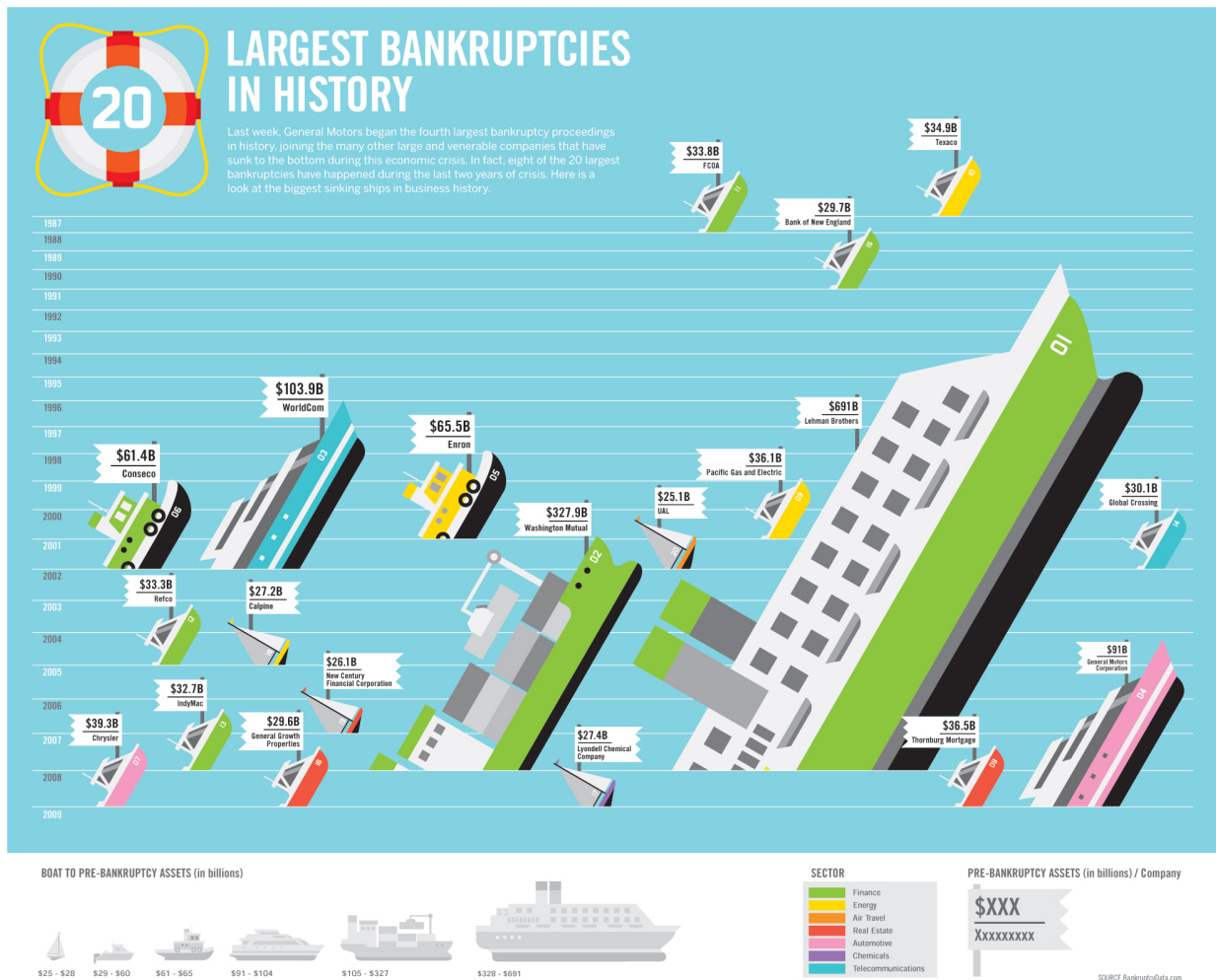
Solution: E. side-by-side boxplots. This allows us to see the distributions of SAT scores per gender and year.

35. [4 Pts.] Consider the following empirical distribution:



- (1) [1 Pt.] The distribution has _____ mode(s).
A. 1 B. 2 **C. 3** D. 4
- (2) [1 Pt.] The distribution is:
A. Skewed left
B. Symmetric
C. Skewed right
- (3) [2 Pts.] Select **all** of the following properties displayed by the distribution:
A. gaps
B. outliers
C. normal left tail
D. None of the above

36. [4 Pts.] Select all of the problems associated with the following plot (there may be more than one problem):



- A. Over-plotting
- B. Use of chart junk**
- C. Vertical axis should be in log scale
- D. Missing vertical axis label**
- E. Poor use of the horizontal dimension**
- F. Graph elements interfere with data**
- G. Stacking
- H. Use of angles to convey information
- I. None of the above are problems with this awesome plot.

End of Exam