

HOMEWORK 2

DECISION TREES, LOGISTIC REGRESSION, LINEAR REGRESSION, AND OPTIMIZATION

CMU 10-701: MACHINE LEARNING (SPRING 2021)

piazza.com/cmu/spring2021/10701

OUT: February 19, 2021

DUE: Thursday, March 04, 2021 11:59pm

TAs: Stefani Karp, Jeffrey Huang

START HERE: Instructions

- **Collaboration policy:** Collaboration on solving the homework is allowed, after you have thought about the problems on your own. It is also OK to get clarification (but not solutions) from books or online resources, again after you have thought about the problems on your own. There are two requirements: first, cite your collaborators fully and completely (e.g., “Jane explained to me what is asked in Question 2.1”). Second, write your solution *independently*: close the book and all of your notes, and send collaborators out of the room, so that the solution comes from you only. See the Academic Integrity Section on the course site for more information:
https://www.cs.cmu.edu/~aarti/Class/10701_Spring21/index.html.
- **Extension Policy:** See the homework extension policy here:
https://www.cs.cmu.edu/~aarti/Class/10701_Spring21/index.html.
- **Submitting your work:**
 - All portions of the assignments should be submitted to Gradescope (<https://gradescope.com/>).
 - **Programming:** We will autograde your Python code in Gradescope. After uploading your code, our grading scripts will autograde your assignment by running your program on a virtual machine (VM). We recommend debugging your implementation on your local machine (or the linux servers) and making sure your code is running correctly before any submission. **Our autograder requires that you write your code using Python 3.6.9 and Numpy 1.17.0.**
 - **Written questions:** **You must type the answers in the provided .tex file. Hand-written solutions will not be accepted. Make sure to answer each question in the provided box. DO NOT change the size of the boxes, because it may mess up the autograder.** Upon submission, make sure to label each question using the template provided by Gradescope. Please make sure to assign ALL pages corresponding to each question.

1 Warm-Up [16 pts]

1. [2 points] For the following decision boundary (where points that satisfy the inequality are the positive class and the rest are the negative class), select all types of models that can perform perfect classification according to this boundary. In the problem below, all points in \mathbb{R}^2 are under consideration.

$$|x|, |y| \leq 1$$

- ☐ K-NN
 - ☐ Decision Tree
 - ☐ Logistic Regression
 - ☐ Linear Regression
2. [2 points] **True or False:** Naive Bayes assumes that all features x_i are independent of all other features.
- ☐ True
 - ☐ False
3. [2 points] **True or False:** In unregularized logistic regression on separable data, checking for 0 change in the weights is always a good stopping criterion.
- ☐ True
 - ☐ False
4. [4 points] **Short Answer:** Under what condition will Gaussian Naive Bayes be guaranteed to find a linear decision boundary for any training set?

5. [6 points] We are given a set of 8 points

$$(1, 1), (2, 4), (3, 9), (4, 16), (5, 25), (6, 34), (7, 43), (8, w)$$

We find that the line $y = mx + h$ best fits this set of 8 points with respect to mean squared loss, for some real numbers m and h , and it has been determined that $m = 4$. What is the value of w ? Show your work.

2 Information Gain & Decision Trees [9 pts]

1. Consider the following data with the following features:

A	B	Y
0	0	1
0	1	0
0	2	0
1	0	1
1	1	0
1	2	1

For all problems in this section, do not show your work.

- (a) **[2 points]** What is the entropy of Y , $H(Y)$, in bits? Note that the unit "bits" implies that you need to use log base 2 in your calculations. (Please round your answer to the third decimal place, e.g. 0.123)

- (b) **[2 points]** Compute the expected information gain of Y upon learning A , $I(Y; A)$. (Please round your answer to the third decimal place, e.g. 0.123)

- (c) **[2 points]** Compute the expected information gain of Y upon learning B , $I(Y; B)$. (Please round your answer to the third decimal place, e.g. 0.123)

- (d) [**3 points**] Apply the greedy decision tree algorithm using information gain as the splitting criterion until the above data are perfectly classified and draw the tree below. Assume that each split may only be based on one feature from the table of data above. Label the non-leaf nodes with which feature the tree will split on, the edges with the value of the attribute, and the leaf nodes with the classification decision. If two different variables give the same information gain, pick the variable which comes first in the alphabet. Your final tree should not perform any unnecessary splits (ones where nothing is gained).

3 Linear Regression [28 points]

In this section, we'll be solving some toy (i.e., small synthetic dataset) linear regression problems, in both closed-form and via iterative optimization (gradient descent, specifically).

Consider $X \in \mathbb{R}^{n \times d}$ (for small values of n and d to be defined below), $y \in \mathbb{R}^n$, and $w \in \mathbb{R}^d$. For simplicity, we will not use a bias term in this problem, but please note that it's typically a very good idea to have a bias term in practice. The objective function we are trying to minimize is:

$$f(w) = \|Xw - y\|_2^2 + \lambda \|w\|_2^2.$$

This is a just a *vectorized* way of computing the square loss over an n -item dataset.

Some of the following questions require **only writing**, and *some* require **writing + code** (each sub-question below specifies exactly which type it is).

We are providing you with `linear_regression.py` as a template, which includes a skeleton gradient descent implementation and incomplete functions corresponding to Parts 3-5 below (you should put your code for Parts 3-5 in the template functions named accordingly; do **not** add your code to this pdf). We will not be running automated tests, but you will be required to submit your completed `linear_regression.py` to Gradescope (for academic integrity purposes).

The only package you need/are allowed to use is `numpy`; for consistency among submitted solutions, please run your code using Python 3.6.9 and `numpy 1.17.0`.

1. **[7 points, writing only]** Show that a closed-form solution for w is $w = (X^\top X + \lambda I)^{-1} X^\top y$ (assuming the inversion is valid).

2. **[1 point, writing only]** It turns out that we have another closed-form solution for w as well: $w = X^\top (X X^\top + \lambda I)^{-1} y$ (assuming the inversion is valid). This is less commonly seen/used in linear regression, but it will likely become relevant later in the course, so it's good to keep in mind.

What is the shape of the matrix we're inverting here, vs. the matrix we inverted in Part 1?

3. **[9 points]** Consider the following design matrix $X = \begin{bmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{bmatrix}$. In the notation above, $n = 3, d = 2$. Consider $y = \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix}$. In what follows, we will try various ways to solve this linear regression problem.

- (a) **[1 point, writing + code]** Attempt to use the closed-form solution in Part 1 to find w , with $\lambda = 0$. Report w or explain what happened if you ran into a problem.

- (b) **[1 point, writing + code]** Attempt to use the closed-form solution in Part 2 to find w , with $\lambda = 0$. Report w or explain what happened if you ran into a problem.

- (c) **[3 points, writing + code]** Now try (a) and (b) again with ℓ_2 regularization at 3 different values of λ : $1, 10^{-3}, 10^{-5}$.

What values of w do you find for (a)? Round to 4 decimal places.

$\lambda = 1$:

$\lambda = 10^{-3}$:

$\lambda = 10^{-5}$:

What values of w do you find for (b)? Round to 4 decimal places.

$\lambda = 1$:

$\lambda = 10^{-3}$:

$\lambda = 10^{-5}$:

What do you notice about this sequence of w values?

- (d) **[3 points, writing + code]** Now use gradient descent to find w (for $\lambda = 0$). What values of w do you find using a learning rate of 0.01 and 10 iterations, 100 iterations, 1000 iterations, 10000 iterations? Round to 4 decimal places.

10 iterations:

100 iterations:

1000 iterations:

10000 iterations:

What do you notice about this sequence of w values?

- (e) **[1 point, writing only]** In one sentence, compare the effect of ℓ_2 regularization in part (c) with the effect of *early stopping* in gradient descent (i.e., running gradient descent for a limited number of iterations) in part (d).

4. **[2 points]** Consider the following design matrix $X = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{bmatrix}$ and $Y = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$. This trivially has a solution: we have 2 equations in 3 unknowns; it is *overparameterized*. It is nevertheless interesting to see how our different approaches behave on this problem.

- (a) **[1 point, writing + code]** Attempt to use the closed-form solution in Part 1 to find w . Report w or explain what happened if you ran into a problem. How has this changed from 3(a)?

- (b) **[1 point, writing + code]** Attempt to use the closed-form solution in Part 2 to find w . Report w or explain what happened if you ran into a problem. How has this changed from 3(b)?

5. **[9 points]** We're now going to examine a slightly different phenomenon - that of conditioning. We'll focus on gradient descent for the remainder of the problem.

We'll consider 2 different X matrices:

$X_1 = \begin{bmatrix} 100 & 0.001 \\ 200 & 0.001 \\ -200 & 0.0005 \end{bmatrix}$ and

$X_2 = \begin{bmatrix} 100 & 100 \\ 200 & 100 \\ -200 & 100 \end{bmatrix}$.

We'll use the same Y for both: $Y = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$.

Complete each of the following questions for *both* X_1 and X_2 , and add 1 sentence of comparison in each part below your results. Put each answer in the labeled box.

- (a) **[3 points, writing + code]** Run gradient descent for 50 iterations with a learning rate of 0.01 and no regularization ($\lambda = 0$). For both $X1$ and $X2$: **a)** What is your final loss value? **b)** What is your final w ?

$X1$ loss:

$X1$ w :

$X2$ loss:

$X2$ w :

1-sentence comparison:

- (b) **[3 points, writing + code]** Recall that the Hessian of a function is its matrix of 2nd-order partial derivatives. For the linear regression objective function above (with $\lambda = 0$), the Hessian is: $2X^T X$. Use `np.linalg.eig` to compute the eigenvalues of the Hessian for $X1$ (first box) and $X2$ (second box).

$X1$:

$X2$:

1-sentence comparison:

- (c) **[3 points, writing + code]** To guarantee convergence of gradient descent, we need to use a learning rate $\leq 1/\lambda_{\max}$. Try setting the learning rate to $1/(2 \times 10^5)$ and rerunning for 50 iterations. For both $X1$ and $X2$: **a)** What loss and **b)** w do you find?

$X1$ loss:

$X1$ w :

$X2$ loss:

$X2$ w :

The Hessian is called ill-conditioned if the ratio of its largest eigenvalue to its smallest eigenvalue is large. Which, if either, of $X1, X2$ yields an ill-conditioned Hessian, and which of $X1, X2$ seems to converge most slowly with the newer learning rate?

4 Logistic Regression [47 pts]

4.1 Logistic Regression Theory [17 points]

1. [5 points] Recall that the **sigmoid function** is defined as

$$\sigma(x) = \frac{1}{1 + e^{-x}}.$$

This is a really important function to understand well!

- a. Plot $\sigma(x)$ and $1 - \sigma(x)$ on the same plot. Restrict the domain to $x \in [-20, 20]$. Make sure to label which curve is $\sigma(x)$ and which curve is $1 - \sigma(x)$.
- b. What relationship between $\sigma(x)$ and $1 - \sigma(x)$ do you notice? (Beyond the fact that the second is one minus the first :) .)
- c. Use this insight to prove a relationship between $\sigma(x)$ and $\sigma(-x)$. Show your work and write your final answer as $\sigma(-x)$ in terms of $\sigma(x)$.

If you solved it correctly, you'll have noticed a very helpful relationship that comes up all the time in logistic regression. Internalizing this is of value!

2. [5 points] Compute $\sigma'(x)$. What does it simplify to? Write your final answer only in terms of $\sigma(x)$. Simplifying is an important part of this problem, because this simplified form is very helpful to know!

3. [7 points] This question focuses on the following: in logistic regression, we sometimes use the relabeling $\{-1, 1\}$ instead of $\{0, 1\}$. We want you to be comfortable seeing both forms and moving fluidly between them.

When our class labels are in $\{0, 1\}$, we can cleverly express $P(Y = y|X = x; \theta)$ as

$$P(Y = 1|X = x; \theta)^y \cdot P(Y = 0|X = x; \theta)^{1-y}.$$

This is the conditional probability for a *single* example.

It yields the following conditional *negative* log likelihood over an n -item dataset:

$$\ell(\theta) = \sum_{i=1}^n \left(-\log \left[P(Y_i = 1|X_i = x_i; \theta)^{y_i} \cdot P(Y_i = 0|X_i = x_i; \theta)^{1-y_i} \right] \right). \quad (1)$$

(Think about what standard, but important, independence assumptions were needed to express the joint conditional probability over the full n -item dataset as a product of

per- i probabilities. You do not have to write anything down for this, but it's important to think about it and not completely take it for granted.)

a. Simplify the i th term of (1) as much as possible when $y_i = 1$. Recall that, in logistic regression, we model $P(Y_i = 1|X_i = x_i; \theta)$, where $\theta = (w, b)$, as $\sigma(w^T x_i + b)$. Your answer should be in terms of x_i, w, b .

b. Simplify the i th term of (1) as much as possible when $y_i = 0$. Recall that, in logistic regression, we model $P(Y_i = 1|X_i = x_i; \theta)$, where $\theta = (w, b)$, as $\sigma(w^T x_i + b)$. Your answer should be in terms of x_i, w, b .

c. Suppose you mapped all appearances of label 0 to label -1 and retained label 1 as is. In light of (a) and (b), what does (1) simplify to? Maximum simplification is part of the credit here. Your answer should be in terms of x_i, y_i, w, b .

For your understanding, can you see why it might be desirable to relabel as $\{-1, 1\}$? You do not have to write/submit anything here. Just think about it.

d. If you relabel as $\{-1, 1\}$ and write $\ell(\theta)$ as in part (c), have you changed the value of $\ell(\theta)$ at all, compared to (1) with $\{0, 1\}$ labels? Justify.

4.2 Logistic Regression Implementation [30 points]

Implementation instructions.

1-Sentence Overview: You will be training a logistic regression model on the Homework 1 dataset by running gradient descent and then answering the written questions below.

Details: You will fill in the `logistic_regression.py` template and submit your complete file to Gradescope, where we will run your code against a suite of tests. Your grade will be automatically determined from the testing results. Since you get immediate feedback after submitting your code and you are allowed to submit as many different versions as you like (without any penalty), it is easy for you to check your code as you go. Our autograder requires that you write your code using **Python 3.6.9 and Numpy 1.17.0**. Otherwise, when running your program on Gradescope, it may produce a result different from the result produced on your local computer. **Please do not include print statements outside the provided functions, as this may crash the autograder.**

Everything to complete is clearly marked in `logistic_regression.py` with the identifier **##### ADD YOUR CODE HERE**. Extensive function descriptions are also provided in the code. Do not otherwise change the functionality of any of the provided code (though you are welcome to replace it if you can guarantee that the functionality won't change).

In this problem, **20 points will come from Gradescope's automated tests**, and **10 points will come from answering the following written questions** using your code.

[10 points] Using your code.

1. [2.5 points] Run gradient descent for 100 epochs with a learning rate of 0.1 (no explicit regularization). **a)** Plot train loss and test loss on one plot and **b.)** plot train accuracy and test accuracy on *another separate* plot. (PLEASE use the provided `_plot` function in the `logistic_regression.py` template.) **c.)** How does your final test accuracy compare to the Naive Bayes test accuracy in Homework 1?

2. **[2.5 points]** If you did this correctly, you might wonder if training longer could squeeze out a bit more test accuracy. Rerun for 3000 epochs and produce *new* versions of the **a.)** loss plot and **b.)** accuracy plot from the previous part. Our implementation took approx. 3 minutes; if yours is taking considerably longer (on your machine), you probably haven't sufficiently vectorized your code. **c.)** Did your test accuracy improve?

3. [5 points] We're now going to study what happens when we replicate features (specifically, we'll replicate 1 feature many times). We'll compare logistic regression to naive Bayes from Homework 1. **For the naive Bayes components, you are encouraged/expected to use your code from Homework 1; you do not need to resubmit the naive Bayes code.** To examine things in more detail and ensure reasonable runtime/convergence/etc., we're going to restrict our *base* dimension to $d = 10$. In other words, our train data matrix will now be 450×10 and our test data matrix will now be 153×10 . Just select the first 10 features from the vocabulary – ensure that you do *not* accidentally sort or shuffle the vocabulary before selecting these 10 features (for autograder purposes).

(a) [2 points] Baseline values (on this 10-feature dataset)

- i. naive Bayes: report the train accuracy and test accuracy of your learned classifier. Clearly label which is train and which is test. No other work required.

- ii. logistic regression: run gradient descent for 5000 epochs with learning rate 0.1, and report the final train accuracy and test accuracy. Clearly label which is train and which is test. No other work required.

(b) [2 points] Now add 500 repeats of the feature at index 4 (0-indexed). Your train data matrix will be 450×510 and your test data matrix will be 153×510 .

- i. naive Bayes: recompute the train accuracy and test accuracy of your learned classifier. Clearly label which is train and which is test. No other work required.

- ii. logistic regression: run gradient descent for 5000 epochs with learning rate 0.1, and compute the final train accuracy and test accuracy. Clearly label which is train and which is test. No other work required.

- (c) [1 point] What do you notice about naive Bayes vs. logistic regression? Please just write 1 sentence.

Collaboration Questions Please answer the following:

1. Did you receive any help whatsoever from anyone in solving this assignment?

Yes / No.

- If you answered 'yes', give full details: _____
- (e.g. "Jane Doe explained to me what is asked in Question 3.4")

2. Did you give any help whatsoever to anyone in solving this assignment?

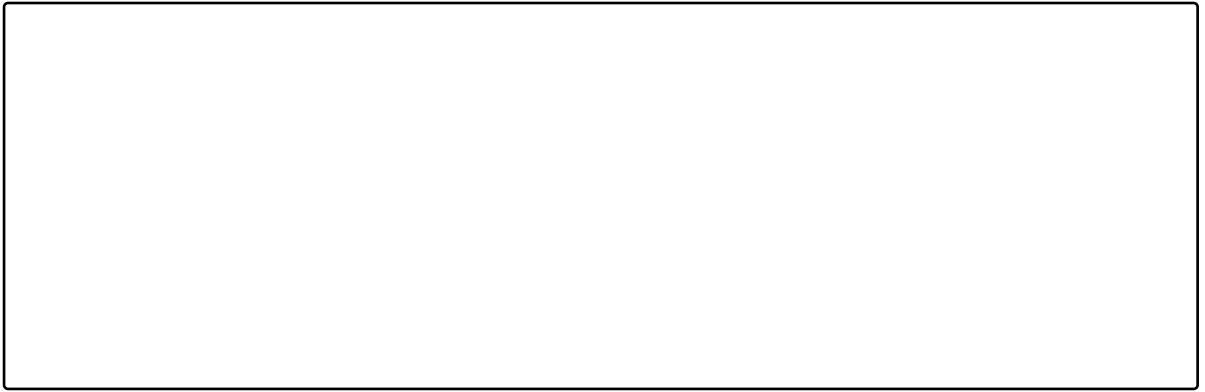
Yes / No.

- If you answered 'yes', give full details: _____
- (e.g. "I pointed Joe Smith to section 2.3 since he didn't know how to proceed with Question 2")

3. Did you find or come across code that implements any part of this assignment ?

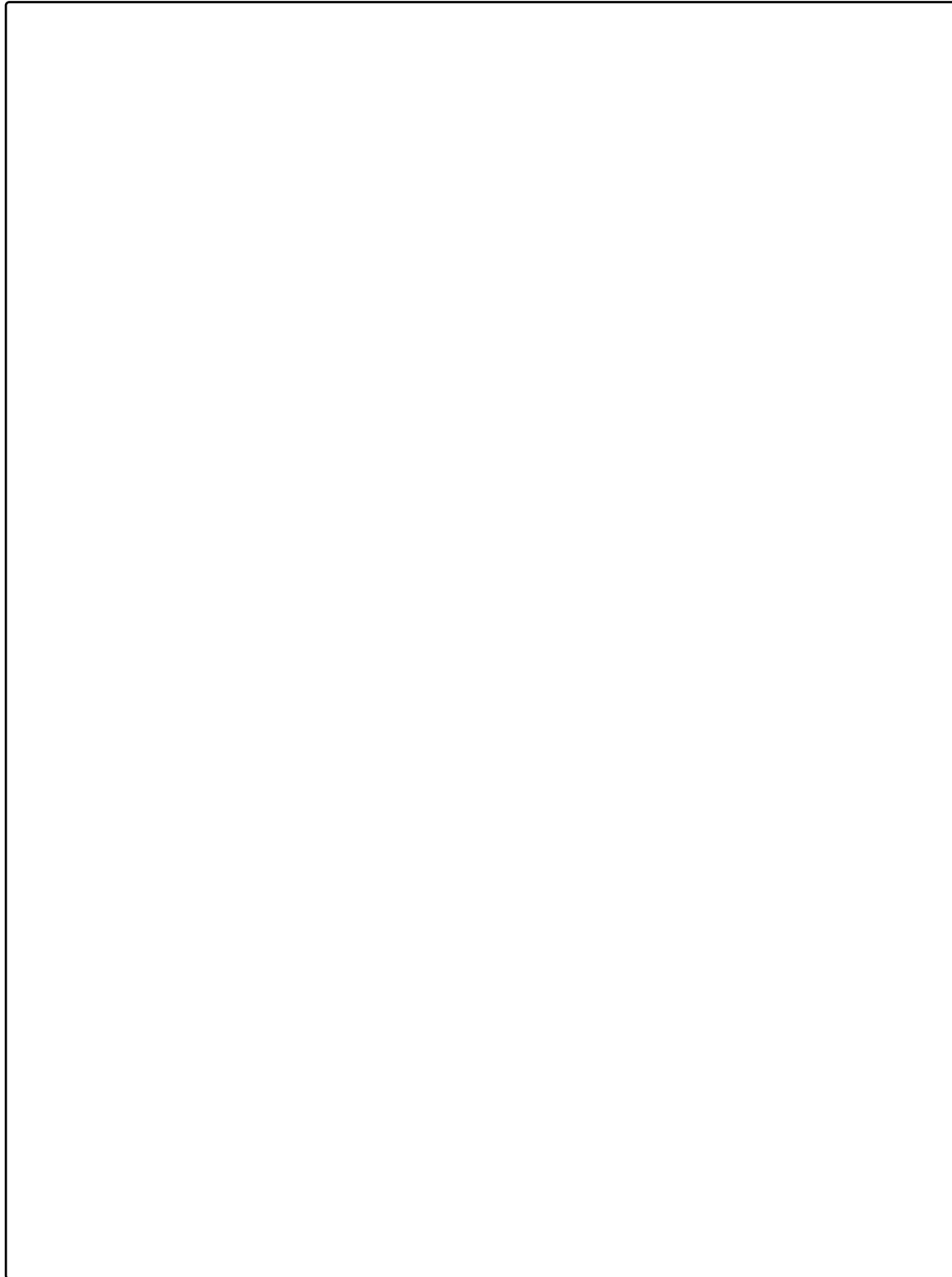
Yes / No. (See below policy on "found code")

- If you answered 'yes', give full details: _____
- (book & page, URL & location within the page, etc.).

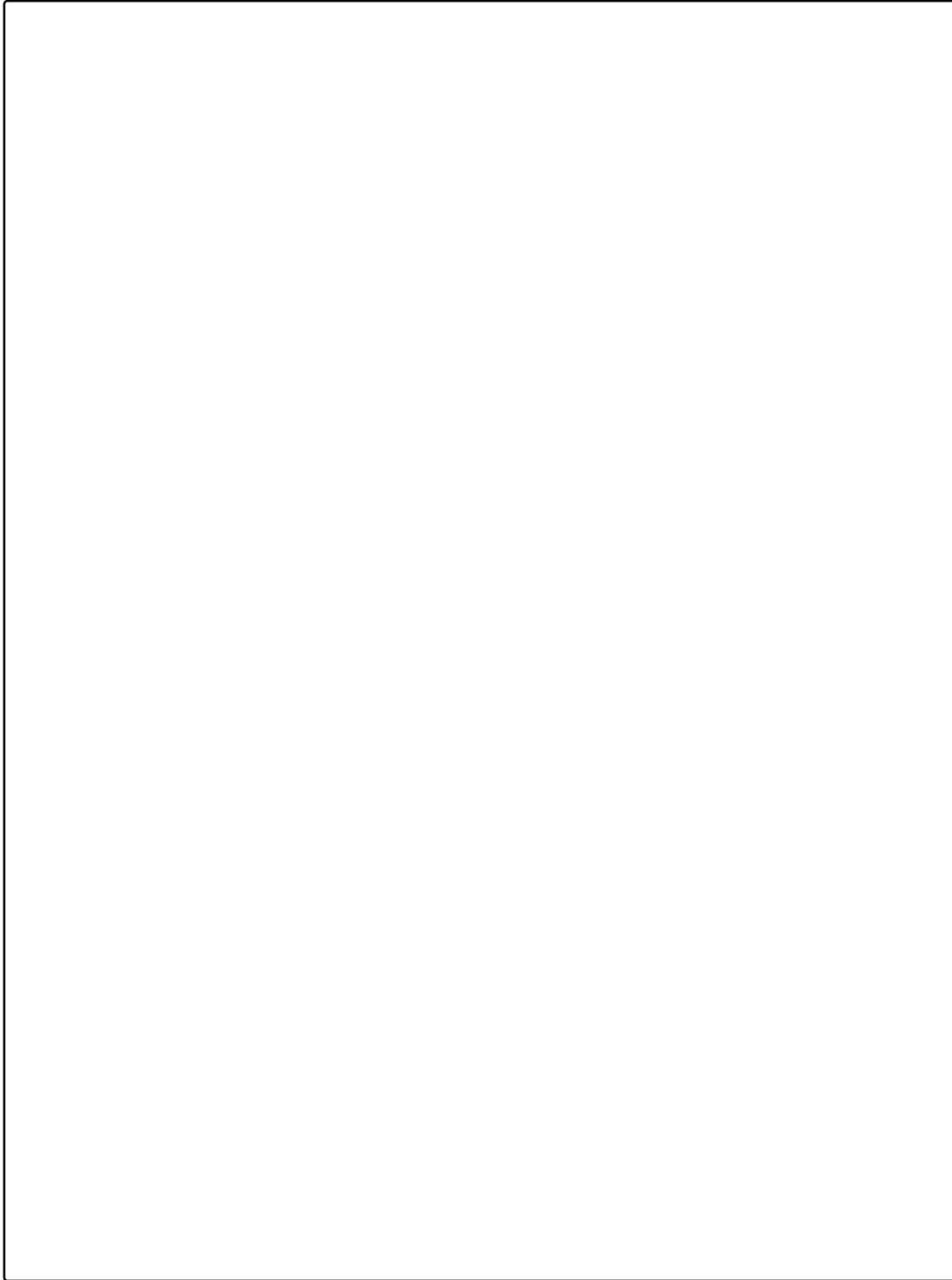


Overflow boxes

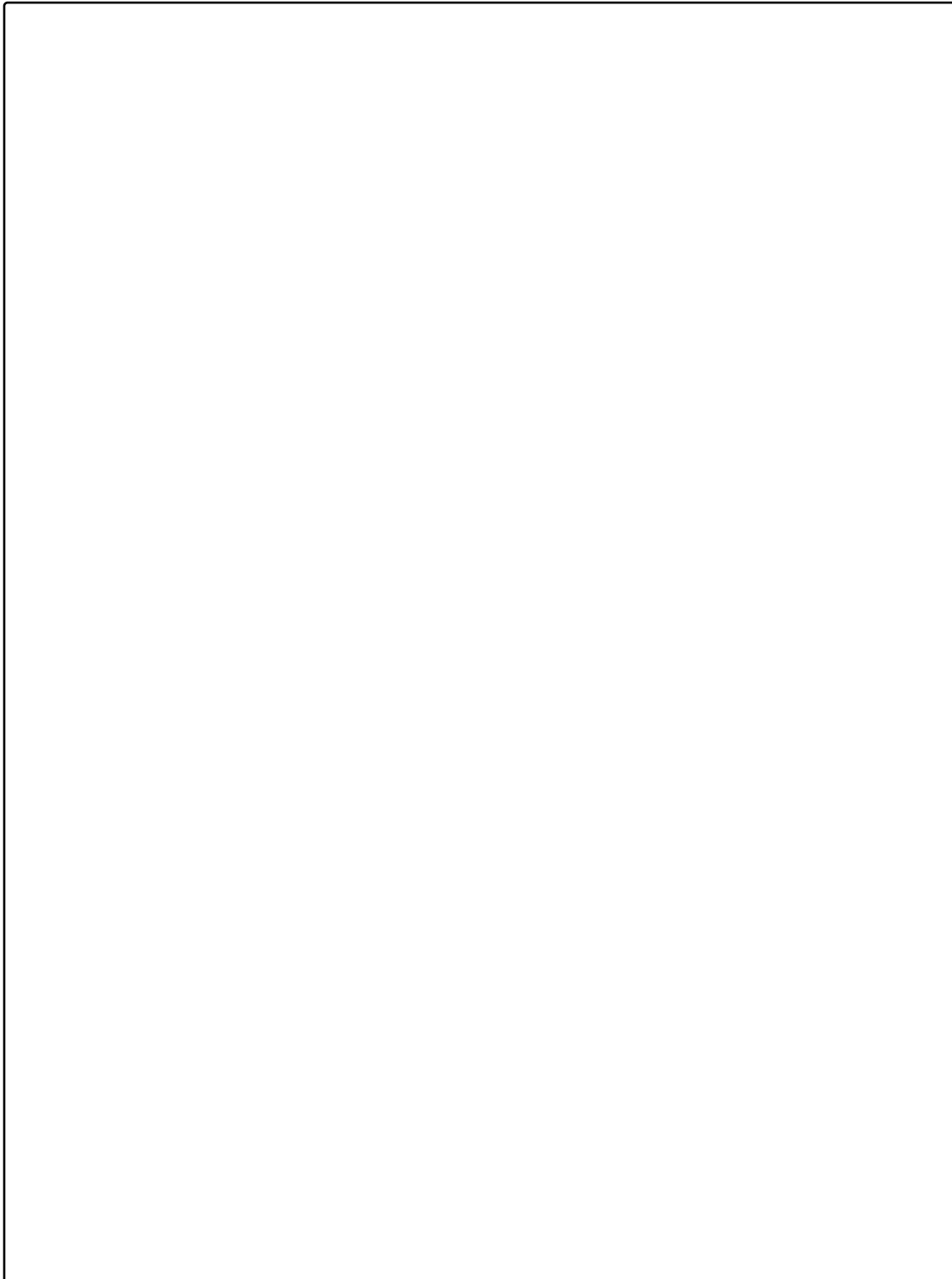
Section 1



Section 2



Section 3



Section 4

